



**Universiteit  
Leiden**  
The Netherlands

## **A transcriptome-wide association study of 229,000 women identifies new candidate susceptibility genes for breast cancer**

Wu, L.; Shi, W.; Long, J.R.; Guo, X.Y.; Michailidou, K.; Beesley, J.; ... ; kConFab AOCS  
Investigators

### **Citation**

Wu, L., Shi, W., Long, J. R., Guo, X. Y., Michailidou, K., Beesley, J., ... Zheng, W. (2018). A transcriptome-wide association study of 229,000 women identifies new candidate susceptibility genes for breast cancer. *Nature Genetics*, 50(7), 968-+. doi:10.1038/s41588-018-0132-x

Version: Not Applicable (or Unknown)  
License: [Leiden University Non-exclusive license](#)  
Downloaded from: <https://hdl.handle.net/1887/86417>

**Note:** To cite this publication please use the final published version (if applicable).



Published in final edited form as:

Nat Genet. 2018 July ; 50(7): 968–978. doi:10.1038/s41588-018-0132-x.

## A transcriptome-wide association study of 229,000 women identifies new candidate susceptibility genes for breast cancer

A full list of authors and affiliations appears at the end of the article.

### Abstract

Breast cancer risk variants identified in genome-wide association studies explain only a small fraction of familial relative risk, and genes responsible for these associations remain largely unknown. To identify novel risk loci and likely causal genes, we performed a transcriptome-wide association study evaluating associations of genetically predicted gene expression with breast

**\*Corresponding Authors:** Wei Zheng, MD, PhD, Division of Epidemiology, Department of Medicine, Vanderbilt Epidemiology Center, Vanderbilt-Ingram Cancer Center, Vanderbilt University Medical Center, 2525 West End Ave, Suite 800, Nashville, Tennessee, 37203, USA. wei.zheng@vanderbilt.edu and Georgia Chenevix-Trench, PhD, Cancer Division, QIMR Berghofer Medical Research Institute, 300 Herston Road, Herston 4006, Australia. Georgia.Trench@qimrberghofer.edu.au.

#### Author Contributions

W.Z. and J.L. conceived the study. L.W. contributed to the study design, and performed statistical analyses. L.W., W.Z. and G.C.-T. wrote the manuscript with significant contributions from W.S., J.L., X.G., and S.L.E.. W.S. performed the *in vitro* experiments. G.C.-T. directed the *in vitro* experiments. X.G. contributed to the model building and pathway analyses. J.B. contributed to the bioinformatics analyses. F.A.-E., E.R., and S.L.E. contributed to the *in vitro* experiments. Y. L. and C. Z. contributed to the model building. K.M., M.K.B., X.-O.S., Q.W., J.D., B.L., C.Z., H.F., A.G., R.T.B., A.M.D., P.D.P.P., J.S., R.L.M., P.K., and D.FE, contributed to manuscript revision, statistical analyses and/or BCAC data management. I.L.A., H.A.-C., V.A., K.J.A., P.L.A., M. Barrdahl, C.B., M.W.B., J.B., M. Bermisheva, C.B., N.V.B., S.E.B., H. Brauch, H. Brenner, L.B., P.B., S.Y.B., B.B., Q.C., T.C., F.C., B.D.C., J.E.C., J.C.-C., X.C., T.-Y.D.C., H.C., C.L.C., NBCS Collaborators, M.C., S.C., F.J.C., D.C., A.C., S.S.C., J.M.C., K.C., M.B.D., P.D., K.F.D., T.D., I.d.S.S., M. Dumont, M. Dwek, D.M.E., U.E., H.E., C.E., M.E., L.F., P.A.F., J.F., D.F.-J., O.F., H.F., L.F., M. Gabrielson, M.G.-D., S.M.G., M.G.-C., M.M.G., M. Ghousaini, G.G.G., M.S.G., D.E.G., A.G.-N., P.G., E. Hahnen, C.A.H., N.H., P. Hall, E. Hallberg, U.H., P. Harrington, A. Hein, B.H., P. Hillemanns, A. Hollestelle, R.N.H., J.L.H., G.H., K.H., D.J.H., A.J., W.J., E.M.J., N.J., K.J., M.E.J., A. Jung, R.K., M.J.K., E.K., V.-M.K., V.N.K., D.L., L.L.M., J. Li, S.L., J. Lissowska, W.-Y.L., S.Loibl, J.L., C.L., M.P.L., R.J.M., T.M., I.M.K., A. Mannermaa, J.E.M., S.M., D.M., H.M.-H., A. Meindl, U.M., J.M., A.M.M., S.L.N., H.N., P.N., S.F.N., B.G.N., O.I.O., J.E.O., H.O., P.P., J.P., D.P.-K., R.P., N.P., K.P., B.R., P.R., N.R., G.R., H.S.R., V.R., A. Romero, J.R., A. Rudolph, E.S., D.P.S., E.J.S., M.K.S., R.K.S., A.S., R.J.S., C. Scott, S.S., M.S., M.J.S., A.S., M.C.S., J.J.S., J.S., H.S., A.J.S., R.T., W.T., J.A.T., M.B.T., D.C.T., A.T., K.T., R.A.E.M.T., D.T., T.T., M.U., C.V., D.V.D.B., D.V., Q.W., C.R.W., C.W., A.S.W., H.W., W.C.W., R.W., A.W., L.X., X.R.Y., A.Z., E.Z., kConFab/AOCS Investigators contributed to the collection of the data and biological samples for the original BCAC studies. All authors have reviewed and approved the final manuscript.

#### URLs.

GTEx protocol, <http://www.gtexportal.org/home/documentationPage>; Gencode V19 annotation file, <http://www.encodegenes.org/releases/19.html>; HaploReg, <http://archive.broadinstitute.org/mammals/haploreg/data/>; OncoArray, <http://epi.grants.cancer.gov/oncoarray/>;

#### Data availability

The GTEx data are publicly available via dbGaP ([www.ncbi.nlm.nih.gov/gap](http://www.ncbi.nlm.nih.gov/gap)); dbGaP Study Accession: phs000424.v6.p1). TCGA data are publicly available via National Cancer Institute's Genomic Data Commons Data Portal (<https://gdc.cancer.gov/>). A subset of the BCAC data that support the findings of this study is publicly available via dbGaP ([www.ncbi.nlm.nih.gov/gap](http://www.ncbi.nlm.nih.gov/gap); accession number phs001265.v1.p1). Most of the BCAC data used in this study are or will be publicly available via dbGaP. Data from some BCAC studies are not publicly available due to restraints imposed by the ethics committees of individual studies; requests for further data can be made to the BCAC (<http://bcac.ccge.medschl.cam.ac.uk/>) Data Access Coordination Committee (DACC). BCAC DACC approval is required to access data from studies ABCFS, ABCS, ABCTB, BBCC, BBBS, BCEES, BCFR-NY, BCFR-PA, BCFR-UT, BCINIS, BSUCH, CBCS, CECILE, CGPS, CTS, DIETCOMPLYF, ESTHER, GC-HBOC, GENICA, GEPARSIXTO, GESBC, HABCS, HCSC, HEBCS, HMBCS, HUBCS, KARBAC, KBBCP, LMBC, MABCS, MARIE, MBCSG, MCBCS, MISS, MMHS, MTLGEBCS, NC-BCFR, OFBCR, ORIGO, pKARMA, POSH, PREFACE, RBCS, SKDKKFZS, SUCCESSB, SUCCESSC, SZBCS, TNBCC, UCIBCS, UKBGS and UKOPS.

#### Code availability

The computer codes used in our study are available upon reasonable request.

#### Competing financial interests

The authors declare no competing financial interests.

cancer risk in 122,977 cases and 105,974 controls of European ancestry. We used data from the Genotype-Tissue Expression Project to establish genetic models to predict gene expression in breast tissue and evaluated model performance using data from The Cancer Genome Atlas. Of the 8,597 genes evaluated, significant associations were identified for 48 at a Bonferroni-corrected threshold of  $P < 5.82 \times 10^{-6}$ , including 14 genes at loci not yet reported for breast cancer. We silenced 13 genes and showed an effect for 11 on cell proliferation and/or colony forming efficiency. Our study provides new insights into breast cancer genetics and biology.

## Keywords

eQTL; genetics; breast cancer; gene expression; GWAS; susceptibility

---

Breast cancer is the most common malignancy among women in many countries<sup>1</sup>. Genetic factors play an important role in its etiology. Since 2007, genome-wide association studies (GWAS) have identified approximately 170 genetic loci harboring common, low-penetrance variants for breast cancer<sup>6-13</sup>, but these variants explain less than 20% of familial relative risk<sup>7</sup>. Most disease-associated risk variants identified by GWAS are located in non-protein coding regions and are not in linkage disequilibrium (LD) with any nonsynonymous coding single nucleotide polymorphisms (SNPs)<sup>14</sup>. Many of these susceptibility variants are located in gene regulatory elements<sup>15,16</sup>, and it has been hypothesized that many GWAS-identified associations may be driven by the regulatory function of risk variants on the expression of nearby genes. For breast cancer, recent studies have already shown that GWAS-identified associations at more than 15 loci are likely due to the effect of risk variants at these loci on regulating the expression of either nearby or more distal genes<sup>7,9,10,13,17-22</sup>. However, for the large majority of the GWAS-identified breast cancer risk loci, the genes responsible for the associations remain unknown.

Several studies have reported that regulatory variants may account for a large proportion of disease heritability not yet discovered through GWAS<sup>23-25</sup>. Many of these variants may have a small effect size, and thus are difficult to identify in individual SNP-based GWAS, even with a large sample size. Applying gene-based approaches that aggregate the effects of multiple variants into a single testing unit may increase study power to identify novel disease-associated loci. Transcriptome-wide association studies (TWAS) systematically investigate the association of genetically predicted gene expression with disease risk, providing an effective approach to identify novel susceptibility genes<sup>26-29</sup>. Recently, Hoffman et al performed a TWAS including 15,440 cases and 31,159 controls and reported significant associations for five genes with breast cancer risk<sup>30</sup>. However, the sample size of that study was relatively small and several reported associations were not significant after Bonferroni correction. Herein, we report results from a larger TWAS of breast cancer that used the MetaXcan method<sup>26</sup> to analyze summary statistics data from 122,977 cases and 105,974 controls of European descent from the Breast Cancer Association Consortium (BCAC).

## Results

### Gene expression prediction models

The study design is shown in Supplementary Figure 1. We used transcriptome and genotyping data from 67 women of European descent included in the Genotype-Tissue Expression (GTEx) project to build genetic models to predict RNA expression levels for each gene expressed in normal breast tissues, by applying the elastic net method ( $\alpha=0.5$ ) with ten-fold cross-validation. Genetically regulated expression was estimated using variants within a 2 MB window flanking the respective gene boundaries, inclusive. SNPs with a minor allele frequency of at least 0.05 and included in the HapMap Phase 2 were used for model building. Of the models built for 12,696 genes, 9,109 showed a prediction performance ( $R^2$ ) of at least 0.01 (10% correlation between predicted and observed expression). For genes for which the expression could not be predicted well using this approach, we built models using only SNPs located in the promoter or enhancer regions, as predicted using three breast cell lines in the Roadmap Epigenomics Project/Encyclopedia of DNA Elements Project. This approach leverages information from functional genomics and reduces the number of variants for variable selection, therefore potentially improving statistical power. This enabled us to build genetic models for additional 3,715 genes with  $R^2 \geq 0.01$ . Supplementary Table 1 provides detailed information regarding the performance threshold and types of models built. Overall, genes that were predicted with  $R^2 \geq 0.01$  in GTEx data were also predicted well in The Cancer Genome Atlas (TCGA) tumor-adjacent normal tissue data (correlation coefficient of 0.55 for  $R^2$  in two datasets; Supplementary Figure 2). Based on model performance in GTEx and TCGA, we prioritized 8,597 genes for analyses of the associations between predicted gene expression and breast cancer risk using the following criteria: 1) genes with a model prediction  $R^2 \geq 0.01$  in the GTEx set (10% correlation) and a Spearman's correlation coefficient of  $\geq 0.1$  in the external validation experiment, 2) genes with a prediction  $R^2 \geq 0.09$  (30% correlation) in the GTEx set regardless of their performance in the TCGA set, 3) genes with a prediction  $R^2 \geq 0.01$  in the GTEx set (10% correlation) that could not be evaluated in the TCGA set because of a lack of data.

### Associations of predicted expression with breast cancer

Using the MetaXcan method<sup>26</sup>, we performed association analyses to evaluate predicted gene expression and breast cancer risk using the meta-analysis summary statistics of SNPs generated for 122,977 cases and 105,974 controls of European ancestry included in BCAC. For the majority of the tested genes, most of the SNPs selected for prediction models were used for the association analyses (e.g., 80% predicting SNPs used for 95.6% of the tested genes). Lambda 1,000 ( $\lambda_{1,000}$ ), a standardized estimate of the genomic inflation scaling to a study of 1,000 cases and 1,000 controls, was 1.004 in our study (Quantile-quantile (QQ) plot presented in Supplementary Figure 3 (a)). Of the 8,597 genes evaluated, we identified 179 whose predicted expression was associated with breast cancer risk at  $P < 1.05 \times 10^{-3}$ , a FDR-corrected significance level (Figure 1, Supplementary Table 2). Of these, 48 showed a significant association at the Bonferroni-corrected threshold of  $P = 5.82 \times 10^{-6}$  (Figure 1, Tables 1–3), including 14 genes located at 11 loci that are 500 kb away from any risk variant identified in previous GWAS (Table 1). An association between lower predicted expression

and increased breast cancer risk was detected for *LRRC3B* (3p24.1), *SPATA18* (4q12), *UBD* (6p22.1), *MIR31HG* (9p21.3), *RIC8A* (11p15.5), *B3GNT1* (11q13.2), *GALNT16* (14q24.1) and *MAN2C1* and *CTD-2323K18.1* (15q24.2). Conversely, an association between higher predicted expression and increased breast cancer risk was identified for *ZSWIM5* (1p34.1), *KLHDC10* (7q32.2), *RP11-867G23.10* (11q13.2), *RP11-218M22.1* (12p13.33) and *PLEKHD1* (14q24.1). The remaining 34 associated genes are located at known breast cancer susceptibility loci (Tables 2–3). Among them, 23 have not yet been implicated as genes responsible for association signals identified at these loci through expression quantitative trait loci (eQTL) and/or functional studies, and do not harbor GWAS or fine-mapping identified risk variants (Table 2), while the other eleven (*KLHDC7A*<sup>7</sup>, *ALS2CR12*<sup>31</sup>, *CASP8*<sup>31,32</sup>, *ATG10*<sup>9</sup>, *SNX3*<sup>233</sup>, *STXBP4*<sup>4,35</sup>, *ZNF404*<sup>8</sup>, *ATP6AP1L*<sup>9</sup>, *RMND117*, *L3MBTL3*<sup>6</sup>, and *RCCD1*<sup>10</sup>) had been reported as potential causal genes at breast cancer susceptibility loci or harbor GWAS or fine-mapping identified risk variants (Table 3). Except for *RP11-7306.3* and *L3MBTL3*, there was no evidence of heterogeneity ( $I^2 < 0.2$ ) across the iCOGS, OncoArray, and GWAS datasets included in our analyses (Supplementary Table 3). Overall, we identified 37 novel susceptibility genes for breast cancer and confirmed eleven genes known to potentially play a role in breast cancer susceptibility.

To determine whether the associations between predicted gene expression and breast cancer risk were independent of GWAS-identified association signals, we performed conditional analyses adjusting for the GWAS-identified risk SNPs closest to the TWAS-identified gene (Supplementary Table 4)<sup>36</sup>. We found that the associations for 11 genes (*LRRC3B*, *SPATA18*, *KLHDC10*, *MIR31HG*, *RIC8A*, *B3GNT1*, *RP11-218M22.1*, *MAN2C1*, *CTD-2323K18.1* (Table 1), *ALK*, *CTD-3051D23.1* (Table 2)) remained statistically significant at  $P < 5.82 \times 10^{-6}$  (Tables 1–3). This suggests the expression of these genes may be associated with breast cancer risk independent of the GWAS-identified risk variant(s). For nine of the genes (*SPATA18*, *KLHDC10*, *MIR31HG*, *RIC8A*, *RP11-218M22.1*, *MAN2C1*, *CTD-2323K18.1* (Table 1), *ALK*, and *CTD-3051D23.1* (Table 2)), the significance of the association remained essentially unchanged, suggesting these associations may be entirely independent of GWAS-identified association signals.

Of the 131 genes showing an association at  $5.82 \times 10^{-6} < P < 1.05 \times 10^{-3}$  (significant after FDR-correction but not Bonferroni-correction), 38 are located at GWAS-identified risk loci (Table 4). Except for *RP11-400F19.8*, there was no evidence of heterogeneity in TWAS association ( $I^2 < 0.2$ ) across the iCOGS, OncoArray, and GWAS studies (Supplementary Table 3). After adjusting for the risk SNPs, associations for *MTHFD1L*, *PVT1*, *RP11-123K19.1*, *FES*, *RP11-400F19.8*, *CTD-2538G9.5*, and *CTD-3216D2.5* remained significant at  $p < 1.05 \times 10^{-3}$ , again suggesting that the association of these genes with breast cancer risk may be independent of the GWAS-identified association signals (Table 4).

For 41 of the 48 associated genes that reached the Bonferroni-corrected significant level, we obtained individual-level data from subjects included in the iCOGS (n=84,740) and OncoArray (n=112,133) datasets, which was 86% of the subjects included in the analysis using summary statistics (Supplementary Table 5). The results from the analysis using individual-level data were very similar to those described above using MetaXcan analyses (Pearson correlation of z-scores was 0.991 for iCOGS data and 0.994 for OncoArray data),

although not all associations reached the Bonferroni-corrected significant level, possibly due to a smaller sample size (Supplementary Table 5). Conditional analyses using individual level data also revealed consistent results compared with analyses using summary data. We found that for several genes within the same genomic region, their predicted expression was correlated with each other (Tables 1–3). The associations between predicted expression of *PLEKHD1* and *ZSWIM5* and breast cancer risk were largely influenced by their corresponding closest risk variants identified in GWAS, although these risk variants are >500 kb away from these genes (Table 1). There were significant correlation of rs999737 and rs1707302 with genetically predicted expression of *PLEKHD1* ( $r = -0.47$  in OncoArray dataset and  $-0.48$  in iCOGS dataset) and *ZSWIM5* ( $r = 0.50$  in OncoArray dataset and  $0.51$  in iCOGS dataset), respectively.

### INQUISIT algorithm scores

For the 48 associated genes after Bonferroni correction, we assessed their integrated expression quantitative trait and *in silico* prediction of GWAS target (INQUISIT) scores<sup>7</sup> to assess whether there are other evidence beyond the scope of eQTL for supporting our TWAS-identified genes as candidate target genes at GWAS-identified loci. The detailed methodology for INQUISIT scores have been described elsewhere<sup>7</sup>. In brief, a score for each gene-SNP pair is calculated across categories representing potential regulatory mechanisms - distal or proximal gene regulation (promoter). Features contributing to the score are based on functionally important genomic annotations such as chromatin interactions, transcription factor binding, and eQTLs. Compared with evidence from eQTL only, INQUISIT scores incorporate additional lines of evidence, including distal regulations. The INQUISIT scores for our identified genes are shown in Supplementary Table 6. Except for *UBD* with a very low score in the distal regulation category (0.05), none of the genes at novel loci (Table 1) showed evidence to be potential target genes for GWAS-identified breast cancer susceptibility loci. This is interesting and within the expectation since these genes may represent novel association signals. There was evidence suggesting that *RP11-439A17.7*, *NUDT17*, *ANKRD34A*, *BTN3A2*, *AP006621.6*, *RPLP2*, *LRRC37A2*, *LRRC37A*, *KANSL1-AS1*, *CRHR1* and *HAPLN4* listed in Table 2, and all eleven genes listed in Table 3, may be target genes for risk variants at these loci (Supplementary Table 6). For *NUDT17*, *ANKRD34A*, *RPLP2*, *LRRC37A2*, *LRRC37A*, *KANSL1-AS1*, *CRHR1*, *HAPLN4*, *KLHDC7A*, *ALS2CR12*, *CASP8*, *ATG10*, *ATP6AP1L*, *L3MBTL3*, *RMND1*, *SNX32*, *RCCD1*, *STXBP4* and *ZNF404*, the INQUISIT scores were not derived only from eQTL data, providing orthogonal support for these genes. For these loci, the associations of candidate causal SNPs with breast cancer risk may be mediated through these genes. This is in general consistent with the findings from the conditional analyses.

### Pathway enrichment analyses

Ingenuity Pathway Analysis (IPA)<sup>37</sup> suggested potential enrichment of cancer-related functions for the identified protein-coding genes (Supplementary Table 7). The top canonical pathways identified included apoptosis related pathways (Granzyme B signaling ( $p=0.024$ ) and cytotoxic T lymphocyte-mediated apoptosis of target cells ( $p=0.046$ )), immune system pathway (inflammasome pathway ( $p=0.030$ )), and tumoricidal function of hepatic natural killer cells ( $p=0.036$ ). The identified pathways are largely consistent with previous findings

<sup>7</sup>. For the associated lncRNAs, pathway analysis of their highly co-expressed protein-coding genes also revealed potential over-representation of cancer-related functions (Supplementary Table 7).

### ***In vitro* assays of gene functions**

To assess the function of genes whose high predicted expression were associated with increased breast cancer risk, we selected 13 genes for knockdown experiments in breast cells: *ZSWIM5*, *KLHDC10*, *RP11-218M22.1* and *PLEKHD1* (Table 1), *UBLCP1*, *AP006621.6*, *RP11-467J12.4*, *CTD-3032H12.1* and *RP11-15A1.7* (Table 2), and *ALS2CR12*, *RMND1*, *STXBP4* and *ZNF404* (Table 3). As negative controls, we selected *B2M*, *ARHGDI1* and *ZAP70* using the criteria: 1) 2 MB from any known breast cancer risk locus; 2) not an essential gene in breast cancer<sup>38,39</sup>; and 3) not predicted to be a target gene in INQUISIT. In addition, as positive controls, we included *PIDD1* (Table 4)<sup>7</sup>, *NRBF2*<sup>20</sup> and *ABHD8*<sup>22</sup>, which have been functionally validated as target genes at breast cancer risk loci. We performed quantitative PCR (qPCR) on a panel of three 'normal' mammary epithelial and 15 breast cancer cell lines to analyze their expression levels (Supplementary Figure 4 and Supplementary Table 8). All 19 genes were expressed in the normal mammary epithelial line 184A1<sup>40</sup> and the luminal breast cancer cell lines, MCF7 and T47D, so we used these cell lines for the proliferation assay, and MCF7 for the colony formation assay<sup>41</sup>. We also evaluated *SNX32*, *ALK* and *BTN3A2* by qPCR, but they were not expressed in T47D and MCF7 cells; therefore they were not evaluated further. It was difficult to design siRNAs against *RP11-867G23.1* and *RP11-53O19.1* because they both have multiple transcripts with limited, GC-rich regions in common. We did not include *RPLP2* because it is already known to be an essential gene for breast cancer survival<sup>42</sup>. Knockdown of the 19 tested genes was achieved by small short interfering RNA (siRNA) (Supplementary Table 9) and the knockdown efficiency was calculated in 184A1, MCF7 and T47D for each siRNA pair. Robust knockdown of the gene of interests (GOI) was validated by qPCR with the majority of the siRNAs (Supplementary Figure 5).

To evaluate the survival and proliferation ability of cells following gene interruption, we used an InCuCyte to quantify cell proliferation in real time and quantified the corrected proliferation of cells with knocking down of GOI in comparison to that of cells with non-target control (NTC) siRNA). As expected, knockdown of the three negative control genes (*B2M*, *ARHGDI1* and *ZAP70*) did not significantly change cell proliferation in any of the three cell lines (Figure 2A, Supplementary Figure 6). However, with the exception of *UBLCP1*, *RMND1* and *STXBP4*, knockdown of all other genes (11 TWAS-identified genes along with two known genes, *ABHD8* and *NRBF2*) resulted in significantly decreased cell proliferation in 184A1 normal breast cells, with *KLHDC10*, *PLEKHD1*, *RP11-218M22.1*, *AP006621.6*, *ZNF404*, *RP11-467J12.4*, *CTD-3032H12.1* and *STXBP4* showing a similar effect in one or both cancer cell lines. Down-regulation of three lncRNAs (*RP11-218M22.1*, *RP11-467J12.4* and *CTD-3032H12.1*) resulted in significant reduction in cell proliferation in all three cell lines. We also evaluated the effect of inhibition of these genes on colony forming ability in MCF7 cells. Knockdown of the three negative control genes did not significantly affect colony forming efficiency (CFE). By contrast, knockdown of *PIDD1*, *RP11-15A1.7*, *RP11-218M22.1*, *AP006621.6*, *ZNF404*, *RP11-467J12.4* and

*CTD-3032H12.1* resulted in significantly decreased CFE in MCF7 cells compared to the NTC (Figure 2B, Supplementary Figure 7).

## Discussion

This is the largest study to systematically evaluate associations of genetically predicted gene expression across the human transcriptome with breast cancer risk. We identified 179 genes showing a significant association at the FDR-corrected significance level. Of these, 48 genes showed an association at the Bonferroni-corrected threshold, including 14 at genomic loci that have not previously been implicated for breast cancer risk. Of the 34 genes located at known risk loci, 23 have not previously been shown to be the targets of GWAS-identified risk SNPs at corresponding loci and not harbor any risk SNPs. Our study provides substantial new information to improve the understanding of genetics and etiology for breast cancer.

It is possible that TWAS-identified genes may be associated with breast cancer through their correlation with disease causal genes. To determine the potential functional significance of TWAS-identified genes and provide evidence for causal inference, we knocked down 13 genes for which high predicted levels of expression were associated with an increased breast cancer risk, in one normal and two breast cancer cell lines, and measured the effect on proliferation and CFE. Although there was some variation between cell lines, knockdown of 11 of the 13 genes showed an effect in at least one cell line, particularly on proliferation in 184A1 normal breast cells; the effects were strongest and most consistent for the lncRNAs, *RP11-218M22.1*, *RP11-467J12.4* and *CTD-3032H12.1*. The observation of a more consistent effect in the normal breast cell line compared with the cancer cell lines is not surprising as cancer cell lines have increased capacity to handle gene interference through mutations which enhance cell survival. Rewiring of pathways and compensatory mechanisms is a hallmark of cancer. Knockdown of *PIDD1*, *NRBF2* and *ABHD8*, for which breast cancer risk associated haplotypes have been shown to be associated with increased expression in reporter assays<sup>7,20,22</sup>, affected either proliferation or colony forming efficiency, supporting the results from this study.

Some of the genes with strong functional evidence from our study have been reported to have important roles in carcinogenesis. For example, *RP11-467J12.4* (PR-lncRNA-1) is a p53-regulated lncRNA that modulates gene expression in response to DNA damage downstream of p53<sup>43</sup>. *STXBP4* encodes Syntaxin binding protein 4, a scaffold protein that can stabilise and prevent degradation of an isoform of p63, a member of the p53 tumor suppressor family<sup>44</sup>. *KLHDC10* encodes a member of the Kelch superfamily that can activate apoptosis signal-regulating kinase 1, contributing to oxidative stress-induced cell death<sup>45</sup>. Notably, another member of this superfamily, *KLHDC7A*, has recently been identified as the target gene at the 1p36 breast cancer risk locus<sup>7</sup>.

*SNX32*, *ALK* and *BTN3A2* are also likely susceptibility genes for breast cancer risk. However, their low or absent expression in our chosen breast cell lines prevented further functional analysis. *ALK* (Anaplastic lymphoma kinase) copy number gain and overexpression have been reported in aggressive and metastatic breast cancers<sup>46</sup>.



Therapeutic targeting of ALK rearrangement has significantly improved survival in advanced ALK-positive lung cancer<sup>47</sup>, making it an attractive target for breast and other cancers. *BTN3A2* is a member of the B7/butyrophilin-like group of Ig superfamily receptors modulating the function of T-lymphocytes. Over-expression of *BTN3A2* in epithelial ovarian cancer is associated with higher infiltrating immune cells and a better prognosis<sup>48</sup>.

Our analyses identified multiple genes with reduced expression associated with increased breast cancer risk. Among them, *LRRC3B* and *CASP8* are putative tumor suppressors in multiple cancers, including breast cancer. Leucine-rich repeat-containing 3B (*LRRC3B*) is a putative LRR-containing transmembrane protein, which is frequently inactivated via promoter hypermethylation leading to inhibition of cancer cell growth, proliferation, and invasion<sup>49</sup>. *CASP8* encodes a member of the cysteine-aspartic acid protease family, which play a central role in cell apoptosis. Previous studies have suggested that caspase-8 may act as a tumor suppressor in certain types of lung cancer and neuroblastoma, although this function has not yet been demonstrated in breast cancer. Notably, several large association studies have identified SNPs at the 2q33/*CASP8* locus associated with increased breast cancer risk<sup>31,50</sup>. Consistent with our data, eQTL analyses showed that the risk alleles for breast cancer were associated with reduced *CASP8* mRNA levels in both peripheral blood lymphocytes and normal breast tissue<sup>31</sup>.

For seven of the genes listed in Tables 1 and 2, we found some evidence from studies using tumor tissues, *in vitro* or *in vivo* experiments linking them to cancer risk (Supplementary Table 10), although their association with breast cancer has not been demonstrated in human studies. For five of them, including *LRRC3B*, *SPATA18*, *RIC8A*, *ALK* and *CRHR1*, previous *in vitro* and *in vivo* experiments and human tissue studies showed a consistent direction of the association as demonstrated in our studies. For two other genes (*UBD* and *MIR31HG*), however, results from previous studies were inconsistent, reporting both potential promoting and inhibiting effects on breast cancer development. Future studies are needed to evaluate functions of these genes.

We included a large number of cases and controls, providing strong statistical power for the association analysis. This large sample size enabled us to identify a large number of candidate breast cancer susceptibility genes, much larger than the number identified in a TWAS study with a sample size of about 20% of ours<sup>30</sup>. The previous study included subjects of different races, which could affect the results as linkage disequilibrium (LD) patterns differ by races. Of the five genes reported in that smaller TWAS that showed a suggestive association with breast cancer risk, the association for the *RCCD1* gene was replicated in our study (Table 3). The other four genes (*ANKLE1*, *DHODH*, *ACAP1* and *LRRC25*) were not evaluated in our study because of unsatisfactory performance of our breast specific models for these genes which were built using the GTEx reference dataset including only female European descendants.

A substantial proportion of SNPs included in the OncoArray and iCOGS were selected from breast cancer GWAS and fine-mapping analyses, and thus these arrays were enriched for association signals with breast cancer risk. As a result, the overall  $\lambda$  value for the BCAC association analyses of individual variants is 1.26 after adjusting for population

stratifications (QQ plot in Supplementary Figure 3 (b))<sup>7</sup>. The  $\lambda$  value for the associations of the ~257,000 SNPs included in the gene expression prediction models of the 8,597 genes tested in our association analysis is 1.40 (QQ plot in Supplementary Figure 3 (c)). This higher  $\lambda$  value is perhaps expected because of a potential further enrichment of breast cancer associated signals in the set of SNPs selected to predict gene expression. There could be additional gain of power (and thus a higher  $\lambda$  value) in TWAS as it aggregates the effect of multiple SNPs to predict gene expression and use genes as the unit for association analyses. The lambda ( $\lambda$ ) for our associated analyses of 8,597 genes was 1.51 (QQ plot presented in Supplementary Figure 3 (a)) likely due to the potential enrichment and power gain as well as our large sample size, and the highly polygenic nature of the disease<sup>7,51</sup>. Interestingly, high  $\lambda$  values were also found in recent large studies of other polygenic traits, such as body mass index (BMI) ( $\lambda = 1.99$ ) and height ( $\lambda = 2.7$ )<sup>52,53</sup>. The  $\lambda_{1,000}$ , a standardized estimate of the genomic inflation scaling to a study of 1,000 cases and 1,000 controls, is 1.004 in our study.

The statistical power of our study is very high to detect associations for genes with a relatively high cis-heritability ( $h^2$ ) (Supplementary Figure 8). For example, our study has 80% statistical power to detect an association with breast cancer risk at  $P < 5.82 \times 10^{-6}$  with an OR of 1.07 or higher per one standard deviation increase (or decrease) in the expression level of genes with an  $h^2$  of 0.1 or higher. One limitation of our study is the small sample size for building gene expression prediction models, which may have affected the precision of model parameter estimates. We expect that models built with a larger sample size will identify additional association signals. We used samples from women of European origin in model building, given differences in gene expression patterns between males and females and in genetic architecture across ethnicities<sup>54</sup>. We also used gene expression data of tumor-adjacent normal tissue samples from European descendants in TCGA as an external validation step to prioritize genes for association analyses. Given potential somatic alterations in tumor-adjacent normal tissues, we retained all models showing a prediction  $R^2$  of at least 0.09 in GTEx, regardless of their performance in TCGA. Not all genes have a significant hereditary component in expression regulation, and thus these genes could not be investigated in our study. For example, previous studies have provided strong evidence to support a significant role of the *TERT*, *ESR1*, *CCND1*, *IGFBP5*, *TET2* and *MRPS30* genes in the etiology of breast cancer. However, expression of these genes cannot be predicted well using the data from female European descendants included in the GTEx and thus they were not included in our association analyses. Supplementary Table 11 summarizes the performance of prediction models and association results for breast cancer target genes reported previously at GWAS-identified loci.

In summary, our study has identified multiple gene candidates that can be further functionally characterized. The silencing experiments we performed suggest that many of the genes identified are likely to mediate risk of breast cancer by affecting proliferation or CFE, two hallmarks of cancer. Further investigation of genes identified in our study will provide additional insight into the biology and genetics of breast cancer.

## Methods

The key elements of the study design, statistical parameters, materials and reagents, and human subjects are included in the **Life Sciences Reporting Summary**.

### Building of gene expression prediction models

We used transcriptome and high-density genotyping data from the Genotype-Tissue Expression (GTEx) study to establish prediction models for genes expressed in normal breast tissues. Details of the GTEx have been described elsewhere<sup>55</sup>. Genomic DNA samples obtained from study subjects included in the GTEx were genotyped using Illumina OMNI 5M or 2.5M SNP Array and RNA samples from 51 tissue sites were sequenced to generate transcriptome profiling data. Genotype data were processed according to the GTEx protocol (see URLs). SNPs with a call rate < 98%, with differential missingness between the two array experiments (5M/2.5M Arrays), with Hardy-Weinberg equilibrium  $p$ -value < 10<sup>-6</sup> (among subjects of European ancestry), or showing batch effects were excluded. One Klinefelter individual, three related individuals, and a chromosome 17 trisomy individual were also excluded. The genotype data were imputed to the Haplotype Reference Consortium reference panel<sup>56</sup> using Minimac3 for imputation and SHAPEIT for prephasing<sup>57,58</sup>. SNPs with high imputation quality ( $r^2 > 0.8$ ), minor allele frequency (MAF) > 0.05, and included in the HapMap Phase 2 version, were used to build expression prediction models. For gene expression data, we used Reads Per Kilobase per Million (RPKM) units from RNA-SeQC<sup>59</sup>. Genes with a median expression level of 0 RPKM across samples were removed, and the RPKM values of each gene were log<sub>2</sub> transformed. We performed quantile normalization to bring the expression profile of each sample to the same scale, and performed inverse quantile normalization for each gene to map each set of expression values to a standard normal. We adjusted for the top ten principal components (PCs) derived from genotype data and the top 15 probabilistic estimation of expression residuals (PEER) factors to correct for batch effects and experimental confounders in model building<sup>60</sup>. Genetic and transcriptome data from 67 female subjects of European descent without a prior breast cancer diagnosis were used to build gene expression prediction models for this study.

We built an expression prediction model for each gene by using the elastic net method as implemented in the glmnet R package, with  $\alpha=0.5$ , as recommended by Gamazon et al<sup>27</sup>. The genetically regulated expression for each gene was estimated by including variants within a 2 MB window flanking the respective gene boundaries, inclusive. Expression prediction models were built for protein coding genes, long non-coding RNAs (lncRNAs), microRNAs (miRNAs), processed transcripts, immunoglobulin genes, and T cell receptor genes, according to categories described in the Gencode V19 annotation file (see URLs). Pseudogenes were not included in the present study because of potential concerns of inaccurate calling<sup>61</sup>. Ten-fold cross-validation was used to validate the models internally. Prediction R<sup>2</sup> values (the square of the correlation between predicted and observed expression) were generated to estimate the prediction performance of each of the gene prediction models established.

For genes that cannot be predicted well using the above approach, we built models using only SNPs located in predicted promoter or enhancer regions in breast cell lines. This approach reduces the number of variants for model building, and thus potentially improves model accuracy, by increasing the ratio of sample size to effective degrees of freedom.

SNP-level annotation data in three breast cell lines, namely, Breast Myoepithelial Primary Cells (E027), Breast variant Human Mammary Epithelial Cells (vHMEC) (E028), and HMEC Mammary Epithelial Primary Cells (E119) in the Roadmap Epigenomics Project/ Encyclopedia of DNA Elements Project<sup>16</sup>, were downloaded from HaploReg (Version 4.0, assessed on December 6, 2016) (see URLs). SNPs in regions classified as promoters (TssA, TssAFlnk), enhancers (Enh, EnhG), or regions with both promoter and enhancer signatures (ExFlnk) according to the core 15 chromatin state model<sup>16</sup> in at least one of the cell lines were retained as input SNPs for model building.

### **Evaluating performance of gene expression prediction models using The Cancer Genome Atlas (TCGA) data**

To assess further the validity of the models, we performed external validation using data generated in tumor-adjacent normal breast tissue samples obtained from 86 European-ancestry female breast cancer patients included in the TCGA. Genotype data were imputed using the same approach as described for GTEx data. Expression data were processed and normalized using a similar approach as described above. The predicted expression level for each gene was calculated using the model established using GTEx data and then compared with the observed level of that gene using the Spearman's correlation.

### **Evaluating statistical power for association tests**

We conducted a simulation analysis to assess the power of our TWAS analysis. Specifically, we set the number of cases and controls to be 122,977 and 105,974, respectively, and generated the gene expression levels from the empirical distribution of predicted gene expression levels in the BCAC. We calculated statistical power at  $P < 5.82 \times 10^{-6}$  (the significance level used in our TWAS) according to cis-heritability ( $h^2$ ) which we aim to capture using gene expression prediction models ( $R^2$ ). The results based on 1000 replicates are summarized in Supplementary Figure 8. Based on the power calculation, our TWAS analysis has 80% power to detect a minimum odds ratio of 1.11, 1.07, 1.05, 1.04, or 1.03 for breast cancer risk per one standard deviation increase (or decrease) in the expression level of a gene whose cis-heritability is 5%, 10%, 20%, 40%, or 60%, respectively.

### **Association analyses of predicted gene expression with breast cancer risk**

We used the following criteria to select genes for the association analysis: 1) with a model prediction  $R^2$  of  $\geq 0.01$  in GTEx and a Spearman's correlation coefficient of  $\geq 0.1$  in TCGA, 2) with a prediction  $R^2$  of  $\geq 0.09$  in GTEx regardless of the performance in TCGA, 3) with a prediction  $R^2$  of  $\geq 0.01$  in GTEx but unable to be evaluated in TCGA. The second group of genes was selected because some gene expression levels might have changed in TCGA tumor-adjacent normal tissues, and thus it is anticipated that some genes may show low prediction performance in TCGA data due to the influence of tumor growth<sup>62,63</sup>. Overall, a

total of 8,597 genes met the criteria and were evaluated for their expression-trait associations.

To identify novel breast cancer susceptibility loci and genes, the MetaXcan method, as described elsewhere, was used for the association analyses<sup>26</sup>. Briefly, the formula:

$$Z_g \approx \sum_{l \in \text{Model}_g} w_{lg} \frac{\hat{\sigma}_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)}$$

was used to estimate the Z-score of the association between predicted expression and breast cancer risk. Here  $w_{lg}$  is the weight of SNP  $l$  for predicting the expression of gene  $g$ ,  $\hat{\beta}_l$  and  $\text{se}(\hat{\beta}_l)$  are the GWAS association regression coefficient and its standard error for SNP  $l$ , and  $\hat{\sigma}_l$  and  $\hat{\sigma}_g$  are the estimated variances of SNP  $l$  and the predicted expression of gene  $g$  respectively. Therefore, the weights for predicting gene expression, GWAS summary statistics results, and correlations between model predicting SNPs are the input variables for the MetaXcan analyses. For this study we estimated correlations between SNPs included in the prediction models using the phase 3, 1000 Genomes Project data focusing on European population.

For the association analysis, we used the summary statistics data of genetic variants associated with breast cancer risk generated in 122,977 breast cancer patients and 105,974 controls of European ancestry from the Breast Cancer Association Consortium (BCAC). The details of the BCAC have been described elsewhere<sup>7,9,13,64,65</sup>. Briefly, 46,785 breast cancer cases and 42,892 controls of European ancestry were genotyped using a custom Illumina iSelect genotyping array (iCOGS) containing ~211,155 variants. A further 61,282 cases and 45,494 controls of European ancestry were genotyped using the OncoArray including 570,000 SNPs (see URLs). Also included in this analysis were data from nine GWAS studies including 14,910 breast cancer cases and 17,588 controls of European ancestry. Genotype data from iCOGS, OncoArray and GWAS were imputed using the October 2014 release of the 1000 Genomes Project data as reference. Genetic association results for breast cancer risk were combined using inverse variance fixed effect meta-analyses<sup>7</sup>. For our study, only SNPs with imputation  $r^2 \geq 0.3$  were used. All participating BCAC studies were approved by their appropriate ethics review boards. Relevant ethical regulations had been complied. This study was approved by the BCAC Data Access Coordination Committee.

Lambda 1,000 ( $\lambda_{1,000}$ ) was calculated to represent a standardized estimate of the genomic inflation scaling to a study of 1,000 cases and 1,000 controls, using the following formula:  $\lambda_{1,000} = 1 + (\lambda_{\text{obs}} - 1) \times (1/n_{\text{cases}} + 1/n_{\text{controls}}) / (1/1,000_{\text{cases}} + 1/1,000_{\text{controls}})$ <sup>66,67</sup>. We used a Bonferroni corrected  $p$  threshold of  $5.82 \times 10^{-6}$  ( $0.05/8,597$ ) to determine a statistically significant association for the primary analyses. To identify additional gene candidates at previously identified susceptibility loci, we also used a false discovery rate (FDR) corrected  $p$  threshold of  $1.05 \times 10^{-3}$  ( $\text{FDR} = 0.05$ ) to determine a significant association. Associated genes with an expression of  $>0.1$  RPKM in less than 10 individuals in GTEx data were excluded as the corresponding prediction models may not be stable.

To determine whether the predicted expression-trait associations were independent of the top signals identified in previous GWAS, we performed GCTA-COJO analyses developed by Yang et al<sup>36</sup> to calculate association betas and standard errors of variants with breast cancer risk after adjusting for the index SNPs of interest. We then re-ran the MetaXcan analyses using the association statistics after conditioning on the index SNPs. This information was used to determine whether the detected expression-trait associations remained significant after adjusting for the index SNPs.

For 41 identified associated genes at the Bonferroni-corrected threshold, we also performed analyses using individual level data in iCOGS (n=84,740) and OncoArray (n=112,133) datasets. We generated predicted gene expression using predicting SNPs (Supplementary Table 12), and then assessed the association between predicted gene expression and breast cancer risk adjusting for study and nine principal components in iCOGS dataset, and country and the first ten principal components in OncoArray dataset. Conditional analyses adjusting for index SNPs were performed to assess potential influence of reported index SNPs on the association between predicted gene expression and breast cancer risk. Furthermore, we evaluated whether the predicted expression levels of genes within a same genomic region were correlated with each other by using the OncoArray data.

### **INQUISIT algorithm scores for TWAS-identified genes**

To evaluate whether there are additional lines of evidence supporting the identified genes as putative target genes of GWAS identified risk SNPs beyond the scope of eQTL, we assessed their INQUISIT algorithm scores, which have been described elsewhere<sup>7</sup>. Briefly, this approach evaluates chromatin interactions between distal and proximal regulatory transcription-factor binding sites and the promoters at the risk regions using Hi-C data generated in HMECs<sup>68</sup> and Chromatin Interaction Analysis by Paired End Tag (ChiA-PET) in MCF7 cells. This could detect genome-wide interactions brought about by, or associated with, CCCTC-binding factor (CTCF), DNA polymerase II (POL2), and Estrogen Receptor (ER), all involved in transcriptional regulation<sup>68</sup>. Annotation of predicted target genes used the Integrated Method for Predicting Enhancer Targets (IM-PET)<sup>69</sup>, the Predicting Specific Tissue Interactions of Genes and Enhancers (PreSTIGE) algorithm<sup>70</sup>, Hnisz<sup>71</sup> and FANTOM<sup>72</sup>. Features contributing to the scores are based on functionally important genomic annotations such as chromatin interactions, transcription factor binding, and eQTLs. The detailed information for the INQUISIT pipeline and scoring strategy has been included in a previous publication<sup>7</sup>. In brief, besides assigning integral points according to different features, we also set up-weighting and down-weighting criteria according to breast cancer driver genes, topologically associated domain (TAD) boundaries, and gene expression levels in relevant breast cell lines. Scores in the distal regulation category range from 0–7, and in the promoter category from 0–4. A score of “none” represents that no evidence was found for regulation of the corresponding gene.

### **Functional enrichment analysis using Ingenuity Pathway Analysis (IPA)**

We performed functional enrichment analysis for the identified protein-coding genes reaching Bonferroni corrected association threshold. To assess potential functionality of the identified lncRNAs, we examined their co-expressed protein-coding genes determined using

expression data of normal breast tissue of European females in GTEx. Spearman's correlations between protein-coding genes and identified lncRNAs of 0.4 or -0.4 were used to indicate a high co-expression. Canonical pathways, top associated diseases and biofunctions, and top networks associated with genes of interest were estimated using IPA software<sup>37</sup>.

### Gene expression in breast cell lines

Total RNA was isolated from 18 cell lines (Supplementary Table 8) using the RNeasy Mini Kit (Qiagen). cDNA was synthesized using the SuperScript III (Invitrogen) and amplified using the Platinum SYBR Green qPCR SuperMix-UDG cocktail (Invitrogen). Two or three primer pairs were used for each gene and the mRNA levels for each sample was measured in technical triplicates for each primer set. The primer sequences are listed in Supplementary Table 13. Experiments were performed using an ABI ViiA(TM) 7 System (Applied Biosystems), and data processing was performed using ABI QuantStudio™ Software V1.1 (Applied Biosystems). The average of Ct from all the primer pairs for each gene was used to calculate  $C_t$ . The relative quantitation of each mRNA normalizing to that in 184A1 was performed using the comparative Ct method ( $2^{-\Delta C_t}$ ) and summarized in Supplementary Figure 4.

### Short interfering RNA (siRNA) silencing

184A1, MCF7 and T47D cells were reverse-transfected with siRNAs targeting genes of interest (GOI) or a non-targeting control siRNA (consi; Shanghai Genepharma) with RNAiMAX (Invitrogen) according to the manufacturer's protocol. Verification of siRNA knockdown of gene expression by qPCR was performed 36 hours after transfection.

### Proliferation and colony formation assays

For proliferation assays, MCF7 and T47D cells were trypsinized at 16 hours post-transfection and seeded into 24 well plates to achieve ~10% confluency. Phase-contrast images were collected with IncuCyte ZOOM (Essen Bioscience) for seven days. Duplicate samples were assessed for each GOI siRNA transfected cells along with non-target control si (NTCsi) treated cells in the same plate. 184A1 cells were reverse-transfected in 96 well plates to achieve 50% confluence at 8 hours after transfection. Two independent experiments were carried out for all siRNAs in all three cell lines. Each cell proliferation time-course was normalized to the baseline confluency and analyzed in GraphPad Prism. The area under the curve was calculated for each concentration (n=4) and used to calculate corrected proliferation (Corrected proliferation % = 100 +/- (relative proliferation in indicated siRNA - proliferation in NTC siRNA) / knockdown efficiency ("+" if the GOI promotes proliferation and "-" if it inhibits proliferation)). For each gene, results from two siRNAs in two independent experiments were averaged and summarized in Figure 2 and Supplementary Figure 6. For colony formation assays; the same number of GOI siRNA transfected MCF7 cells was seeded in 6 well plates at 16 hours after transfection to assay colony forming efficiency at two weeks. All siRNA-treated cells were seeded in duplicate. Colonies (defined to consist of at least 50 cells) were fixed with methanol, stained with crystal violet (0.5% w/v), scanned and counted using ImageJ as batch analysis by a self-defined plug-in Macro. Correct CFE % = 100 +/- (relative CFE in indicated siRNA - CFE in

NTC siRNA) / knockdown efficiency (“+” if the GOI promotes CF and “-” if it inhibits CF). For each gene, results from two siRNAs in two independent experiments were averaged and summarized in Figure 2 and Supplementary Figure 7. *P*-values were determined by one-way ANOVA followed by Dunnett’s multiple comparisons test.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Authors

Lang Wu<sup>1,163</sup>, Wei Shi<sup>2,163</sup>, Jirong Long<sup>1</sup>, Xingyi Guo<sup>1</sup>, Kyriaki Michailidou<sup>3,4</sup>, Jonathan Beesley<sup>2</sup>, Manjeet K. Bolla<sup>3</sup>, Xiao-Ou Shu<sup>1</sup>, Yingchang Lu<sup>1</sup>, Qiuyin Cai<sup>1</sup>, Fares Al-Ejeh<sup>2</sup>, Esdy Rozali<sup>2</sup>, Qin Wang<sup>3</sup>, Joe Dennis<sup>3</sup>, Bingshan Li<sup>151</sup>, Chenjie Zeng<sup>1</sup>, Helian Feng<sup>5,6</sup>, Alexander Gusev<sup>153,154,155</sup>, Richard T. Barfield<sup>5</sup>, Irene L. Andrulis<sup>7,8</sup>, Hoda Anton-Culver<sup>9</sup>, Volker Arndt<sup>10</sup>, Kristan J. Aronson<sup>11</sup>, Paul L. Auer<sup>12,13</sup>, Myrto Barrdahl<sup>14</sup>, Caroline Baynes<sup>15</sup>, Matthias W. Beckmann<sup>16</sup>, Javier Benitez<sup>17,18</sup>, Marina Bermisheva<sup>19,20</sup>, Carl Blomqvist<sup>21,159</sup>, Natalia V. Bogdanova<sup>20,22,23</sup>, Stig E. Bojesen<sup>24,25,26</sup>, Hiltrud Brauch<sup>27,28,29</sup>, Hermann Brenner<sup>10,29,30</sup>, Louise Brinton<sup>31</sup>, Per Broberg<sup>32</sup>, Sara Y. Brucker<sup>33</sup>, Barbara Burwinkel<sup>34,35</sup>, Trinidad Caldés<sup>36</sup>, Federico Canzian<sup>37</sup>, Brian D. Carter<sup>38</sup>, J. Esteban Castelao<sup>39</sup>, Jenny Chang-Claude<sup>14,40</sup>, Xiaoqing Chen<sup>2</sup>, Ting-Yuan David Cheng<sup>41</sup>, Hans Christiansen<sup>22</sup>, Christine L. Clarke<sup>42</sup>, NBCS Collaborators<sup>43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124</sup>, Margriet Collée<sup>46</sup>, Sten Cornelissen<sup>47</sup>, Fergus J. Couch<sup>48</sup>, David Cox<sup>49,50</sup>, Angela Cox<sup>51</sup>, Simon S. Cross<sup>52</sup>, Julie M. Cunningham<sup>48</sup>, Kamila Czene<sup>53</sup>, Mary B. Daly<sup>54</sup>, Peter Devilee<sup>55,56</sup>, Kimberly F. Doherty<sup>57</sup>, Thilo Dörk<sup>20</sup>, Isabel dos-Santos-Silva<sup>58</sup>, Martine Dumont<sup>59</sup>, Miriam Dwek<sup>60</sup>, Diana M. Eccles<sup>61</sup>, Ursula Eilber<sup>14</sup>, A. Heather Eliassen<sup>6,62</sup>, Christoph Engel<sup>63</sup>, Mikael Eriksson<sup>53</sup>, Laura Fachal<sup>15</sup>, Peter A. Fasching<sup>16,64</sup>, Jonine Figueroa<sup>31,65</sup>, Dieter Flesch-Janys<sup>66,67</sup>, Olivia Fletcher<sup>68</sup>, Henrik Flyger<sup>69</sup>, Lin Fritschi<sup>70</sup>, Marike Gabrielson<sup>53</sup>, Manuela Gago-Dominguez<sup>71,72</sup>, Susan M. Gapstur<sup>38</sup>, Montserrat García-Closas<sup>31</sup>, Mia M. Gaudet<sup>38</sup>, Maya Ghossaini<sup>15</sup>, Graham G. Giles<sup>73,74</sup>, Mark S. Goldberg<sup>75,76</sup>, David E. Goldgar<sup>77</sup>, Anna González-Neira<sup>17</sup>, Pascal Guénel<sup>78</sup>, Eric Hahnen<sup>79,80,81</sup>, Christopher A. Haiman<sup>82</sup>, Niclas Håkansson<sup>83</sup>, Per Hall<sup>53,161</sup>, Emily Hallberg<sup>84</sup>, Ute Hamann<sup>85</sup>, Patricia Harrington<sup>15</sup>, Alexander Hein<sup>16</sup>, Belynda Hicks<sup>86</sup>, Peter Hillemanns<sup>20</sup>, Antoinette Hollestelle<sup>87</sup>, Robert N. Hoover<sup>31</sup>, John L. Hopper<sup>74</sup>, Guanmengqian Huang<sup>85</sup>, Keith Humphreys<sup>53</sup>, David J. Hunter<sup>6,158</sup>, Anna Jakubowska<sup>88,162</sup>, Wolfgang Janni<sup>89</sup>, Esther M. John<sup>90,91,92</sup>, Nichola Johnson<sup>68</sup>, Kristine Jones<sup>86</sup>, Michael E. Jones<sup>93</sup>, Audrey Jung<sup>14</sup>, Rudolf Kaaks<sup>14</sup>, Michael J. Kerin<sup>94</sup>, Elza Khusnutdinova<sup>19,95</sup>, Veli-Matti Kosma<sup>96,97,98</sup>, Vessela N. Kristensen<sup>99,100,101</sup>, Diether Lambrechts<sup>102,103</sup>, Loic Le Marchand<sup>104</sup>, Jingmei Li<sup>157</sup>, Sara Lindström<sup>105,160</sup>, Jolanta Lissowska<sup>106</sup>, Wing-Yee Lo<sup>27,28</sup>, Sibylle Loibl<sup>107</sup>, Jan Lubinski<sup>88</sup>, Craig Luccarini<sup>15</sup>, Michael P. Lux<sup>16</sup>, Robert J.



MacInnis<sup>73,74</sup>, Tom Maishman<sup>108</sup>, Ivana Maleva Kostovska<sup>20,109</sup>, Arto Mannermaa<sup>96,97,98</sup>, JoAnn E. Manson<sup>6,110</sup>, Sara Margolin<sup>111</sup>, Dimitrios Mavroudis<sup>112</sup>, Hanne Meijers-Heijboer<sup>152</sup>, Alfons Meindl<sup>113</sup>, Usha Menon<sup>114</sup>, Jeffery Meyer<sup>48</sup>, Anna Marie Mulligan<sup>115,116</sup>, Susan L. Neuhausen<sup>117</sup>, Heli Nevanlinna<sup>118</sup>, Patrick Neven<sup>119</sup>, Sune F. Nielsen<sup>24,25</sup>, Børge G. Nordestgaard<sup>24,25,26</sup>, Olufunmilayo I. Olopade<sup>120</sup>, Janet E. Olson<sup>84</sup>, Håkan Olsson<sup>32</sup>, Paolo Peterlongo<sup>121</sup>, Julian Peto<sup>58</sup>, Dijana Plaseska-Karanfilska<sup>109</sup>, Ross Prentice<sup>12</sup>, Nadege Presneau<sup>60</sup>, Katri Pylkäs<sup>122,123</sup>, Brigitte Rack<sup>89</sup>, Paolo Radice<sup>125</sup>, Nazneen Rahman<sup>126</sup>, Gad Rennert<sup>127</sup>, Hedy S. Rennert<sup>127</sup>, Valerie Rhenius<sup>15</sup>, Atocha Romero<sup>36,128</sup>, Jane Romm<sup>57</sup>, Anja Rudolph<sup>14</sup>, Emmanouil Saloustros<sup>129</sup>, Dale P. Sandler<sup>130</sup>, Elinor J. Sawyer<sup>131</sup>, Marjanka K. Schmidt<sup>47,132</sup>, Rita K. Schmutzler<sup>79,80,81</sup>, Andreas Schneeweiss<sup>34,133</sup>, Rodney J. Scott<sup>134,135</sup>, Christopher G. Scott<sup>84</sup>, Sheila Seal<sup>126</sup>, Mitul Shah<sup>15</sup>, Martha J. Shrubsole<sup>1</sup>, Ann Smeets<sup>119</sup>, Melissa C. Southey<sup>136</sup>, John J. Spinelli<sup>137,138</sup>, Jennifer Stone<sup>139,140</sup>, Harald Surowy<sup>34,35</sup>, Anthony J. Swerdlow<sup>93,141</sup>, Rulla M. Tamimi<sup>5,6,62</sup>, William Tapper<sup>61</sup>, Jack A. Taylor<sup>130,142</sup>, Mary Beth Terry<sup>143</sup>, Daniel C. Tessier<sup>144</sup>, Abigail Thomas<sup>84</sup>, Kathrin Thöne<sup>40</sup>, Rob A.E.M. Tollenaar<sup>145</sup>, Diana Torres<sup>85,146</sup>, Thérèse Truong<sup>78</sup>, Michael Untch<sup>147</sup>, Celine Vachon<sup>84</sup>, David Van Den Berg<sup>82</sup>, Daniel Vincent<sup>144</sup>, Quinten Waisfisz<sup>152</sup>, Clarice R. Weinberg<sup>148</sup>, Camilla Wendt<sup>111</sup>, Alice S. Whittemore<sup>91,92</sup>, Hans Wildiers<sup>119</sup>, Walter C. Willett<sup>6,62,156</sup>, Robert Winqvist<sup>122,123</sup>, Alicja Wolk<sup>83</sup>, Lucy Xia<sup>82</sup>, Xiaohong R. Yang<sup>31</sup>, Argyrios Ziogas<sup>9</sup>, Elad Ziv<sup>149</sup>, kConFab/AOCS Investigators<sup>150</sup>, Alison M. Dunning<sup>15</sup>, Paul D.P. Pharoah<sup>3,15</sup>, Jacques Simard<sup>59</sup>, Roger L. Milne<sup>73,74</sup>, Stacey L. Edwards<sup>2</sup>, Peter Kraft<sup>5,6</sup>, Douglas F. Easton<sup>3,15</sup>, Georgia Chenevix-Trench<sup>2,\*</sup>, and Wei Zheng<sup>1,\*</sup>

## Affiliations

<sup>1</sup>Division of Epidemiology, Department of Medicine, Vanderbilt Epidemiology Center, Vanderbilt-Ingram Cancer Center, Vanderbilt University School of Medicine, Nashville, TN, USA. <sup>2</sup>Cancer Division, QIMR Berghofer Medical Research Institute, Brisbane, Australia. <sup>3</sup>Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK. <sup>4</sup>Department of Electron Microscopy/Molecular Pathology, The Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus. <sup>5</sup>Program in Genetic Epidemiology and Statistical Genetics, Harvard T.H. Chan School of Public Health, Boston, MA, USA. <sup>6</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA. <sup>7</sup>Fred A. Litwin Center for Cancer Genetics, Lunenfeld-Tanenbaum Research Institute of Mount Sinai Hospital, Toronto, ON, Canada. <sup>8</sup>Department of Molecular Genetics, University of Toronto, Toronto, ON, Canada. <sup>9</sup>Department of Epidemiology, University of California Irvine, Irvine, CA, USA. <sup>10</sup>Division of Clinical Epidemiology and Aging Research, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>11</sup>Department of Public Health Sciences, and Cancer Research Institute, Queen's University, Kingston, ON, Canada. <sup>12</sup>Cancer Prevention Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>13</sup>Zilber School of Public Health, University of Wisconsin-Milwaukee, Milwaukee, WI, USA. <sup>14</sup>Division of Cancer Epidemiology, German Cancer Research Center (DKFZ),

Heidelberg, Germany. <sup>15</sup>Centre for Cancer Genetic Epidemiology, Department of Oncology, University of Cambridge, Cambridge, UK. <sup>16</sup>Department of Gynaecology and Obstetrics, University Hospital Erlangen, Friedrich-Alexander University Erlangen-Nuremberg, Comprehensive Cancer Center Erlangen-EMN, Erlangen, Germany. <sup>17</sup>Human Cancer Genetics Program, Spanish National Cancer Research Centre, Madrid, Spain. <sup>18</sup>Centro de Investigación en Red de Enfermedades Raras (CIBERER), Valencia, Spain. <sup>19</sup>Institute of Biochemistry and Genetics, Ufa Scientific Center of Russian Academy of Sciences, Ufa, Russia. <sup>20</sup>Gynaecology Research Unit, Hannover Medical School, Hannover, Germany. <sup>21</sup>Department of Oncology, Helsinki University Hospital, University of Helsinki, Helsinki, Finland. <sup>22</sup>Department of Radiation Oncology, Hannover Medical School, Hannover, Germany. <sup>23</sup>N.N. Alexandrov Research Institute of Oncology and Medical Radiology, Minsk, Belarus. <sup>24</sup>Copenhagen General Population Study, Herlev and Gentofte Hospital, Copenhagen University Hospital, Herlev, Denmark. <sup>25</sup>Department of Clinical Biochemistry, Herlev and Gentofte Hospital, Copenhagen University Hospital, Herlev, Denmark. <sup>26</sup>Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark. <sup>27</sup>Dr. Margarete Fischer-Bosch-Institute of Clinical Pharmacology, Stuttgart, Germany. <sup>28</sup>University of Tübingen, Tübingen, Germany. <sup>29</sup>German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>30</sup>Division of Preventive Oncology, German Cancer Research Center (DKFZ) and National Center for Tumor Diseases (NCT), Heidelberg, Germany. <sup>31</sup>Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA. <sup>32</sup>Department of Cancer Epidemiology, Clinical Sciences, Lund University, Lund, Sweden. <sup>33</sup>Department of Gynecology and Obstetrics, University of Tübingen, Tübingen, Germany. <sup>34</sup>Department of Obstetrics and Gynecology, University of Heidelberg, Heidelberg, Germany. <sup>35</sup>Molecular Epidemiology Group, C080, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>36</sup>Medical Oncology Department, CIBERONC Hospital Clínico San Carlos, Madrid, Spain. <sup>37</sup>Genomic Epidemiology Group, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>38</sup>Epidemiology Research Program, American Cancer Society, Atlanta, GA, USA. <sup>39</sup>Oncology and Genetics Unit, Instituto de Investigación Biomedica Galicia Sur (IISGS), Xerencia de Xestión Integrada de Vigo-SERGAS, Vigo, Spain. <sup>40</sup>University Cancer Center Hamburg (UCCH), University Medical Center Hamburg-Eppendorf, Hamburg, Germany. <sup>41</sup>Department of Epidemiology, University of Florida, Gainesville, FL, USA. <sup>42</sup>Westmead Institute for Medical Research, University of Sydney, Sydney, Australia. <sup>43</sup>Department of Oncology, Haukeland University Hospital, Bergen, Norway. <sup>44</sup>National Advisory Unit on Late Effects after Cancer Treatment, Oslo University Hospital Radiumhospitalet, Oslo, Norway. <sup>45</sup>Oslo University Hospital, Oslo, Norway. <sup>46</sup>Department of Clinical Genetics, Erasmus University Medical Center, Rotterdam, The Netherlands. <sup>47</sup>Division of Molecular Pathology, The Netherlands Cancer Institute - Antoni van Leeuwenhoek Hospital, Amsterdam, The Netherlands. <sup>48</sup>Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester, MN, USA. <sup>49</sup>Department of Epidemiology and Biostatistics,

School of Public Health, Imperial College London, London, UK. <sup>50</sup>.INSERM U1052, Cancer Research Center of Lyon, Lyon, France. <sup>51</sup>.Sheffield Institute for Nucleic Acids, Department of Oncology and Metabolism, University of Sheffield, Sheffield, UK. <sup>52</sup>.Academic Unit of Pathology, Department of Neuroscience, University of Sheffield, Sheffield, UK. <sup>53</sup>.Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden. <sup>54</sup>.Department of Clinical Genetics, Fox Chase Cancer Center, Philadelphia, PA, USA. <sup>55</sup>.Department of Pathology, Leiden University Medical Center, Leiden, The Netherlands. <sup>56</sup>.Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands. <sup>57</sup>.Center for Inherited Disease Research (CIDR), Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA. <sup>58</sup>.Department of Non-Communicable Disease Epidemiology, London School of Hygiene and Tropical Medicine, London, UK. <sup>59</sup>.Genomics Center, Centre Hospitalier Universitaire de Québec Research Center, Laval University, Québec City, QC, Canada. <sup>60</sup>.Department of Biomedical Sciences, Faculty of Science and Technology, University of Westminster, London, UK. <sup>61</sup>.Cancer Sciences Academic Unit, Faculty of Medicine, University of Southampton, Southampton, UK. <sup>62</sup>.Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>63</sup>.Institute for Medical Informatics, Statistics and Epidemiology, University of Leipzig, Leipzig, Germany. <sup>64</sup>.David Geffen School of Medicine, Department of Medicine Division of Hematology and Oncology, University of California at Los Angeles, Los Angeles, CA, USA. <sup>65</sup>.Usher Institute of Population Health Sciences and Informatics, The University of Edinburgh Medical School, Edinburgh, UK. <sup>66</sup>.Institute for Medical Biometrics and Epidemiology, University Medical Center Hamburg-Eppendorf, Hamburg, Germany. <sup>67</sup>.Department of Cancer Epidemiology, Clinical Cancer Registry, University Medical Center Hamburg-Eppendorf, Hamburg, Germany. <sup>68</sup>.The Breast Cancer Now Toby Robins Research Centre, The Institute of Cancer Research, London, UK. <sup>69</sup>.Department of Breast Surgery, Herlev and Gentofte Hospital, Copenhagen University Hospital, Herlev, Denmark. <sup>70</sup>.School of Public Health, Curtin University, Perth, Australia. <sup>71</sup>.Genomic Medicine Group, Galician Foundation of Genomic Medicine, Instituto de Investigación Sanitaria de Santiago de Compostela (IDIS), Complejo Hospitalario Universitario de Santiago, SERGAS, Santiago De Compostela, Spain. <sup>72</sup>.Moore's Cancer Center, University of California San Diego, La Jolla, CA, USA. <sup>73</sup>.Cancer Epidemiology & Intelligence Division, Cancer Council Victoria, Melbourne, Australia. <sup>74</sup>.Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Melbourne, Australia. <sup>75</sup>.Department of Medicine, McGill University, Montréal, QC, Canada. <sup>76</sup>.Division of Clinical Epidemiology, Royal Victoria Hospital, McGill University, Montréal, QC, Canada. <sup>77</sup>.Department of Dermatology, Huntsman Cancer Institute, University of Utah School of Medicine, Salt Lake City, UT, USA. <sup>78</sup>.Cancer & Environment Group, Center for Research in Epidemiology and Population Health (CESP), INSERM, University Paris-Sud, University Paris-Saclay, Villejuif, France. <sup>79</sup>.Center for Hereditary Breast and Ovarian Cancer, University Hospital of Cologne,

Cologne, Germany. <sup>80</sup>Center for Integrated Oncology (CIO), University Hospital of Cologne, Cologne, Germany. <sup>81</sup>Center for Molecular Medicine Cologne (CMMC), University of Cologne, Cologne, Germany. <sup>82</sup>Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>83</sup>Institute of Environmental Medicine, Karolinska Institutet, Stockholm, Sweden. <sup>84</sup>Department of Health Sciences Research, Mayo Clinic, Rochester, MN, USA. <sup>85</sup>Molecular Genetics of Breast Cancer, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>86</sup>Cancer Genomics Research Laboratory, Leidos Biomedical Research, Frederick National Laboratory for Cancer Research, Frederick, MD, USA. <sup>87</sup>Department of Medical Oncology, Family Cancer Clinic, Erasmus MC Cancer Institute, Rotterdam, The Netherlands. <sup>88</sup>Department of Genetics and Pathology, Pomeranian Medical University, Szczecin, Poland. <sup>89</sup>Department of Gynecology and Obstetrics, University Hospital Ulm, Ulm, Germany. <sup>90</sup>Department of Epidemiology, Cancer Prevention Institute of California, Fremont, CA, USA. <sup>91</sup>Department of Health Research and Policy - Epidemiology, Stanford University School of Medicine, Stanford, CA, USA. <sup>92</sup>Stanford Cancer Institute, Stanford University School of Medicine, Stanford, CA, USA. <sup>93</sup>Division of Genetics and Epidemiology, The Institute of Cancer Research, London, UK. <sup>94</sup>School of Medicine, National University of Ireland, Galway, Ireland. <sup>95</sup>Department of Genetics and Fundamental Medicine, Bashkir State University, Ufa, Russia. <sup>96</sup>Translational Cancer Research Area, University of Eastern Finland, Kuopio, Finland. <sup>97</sup>Institute of Clinical Medicine, Pathology and Forensic Medicine, University of Eastern Finland, Kuopio, Finland. <sup>98</sup>Imaging Center, Department of Clinical Pathology, Kuopio University Hospital, Kuopio, Finland. <sup>99</sup>Department of Cancer Genetics, Institute for Cancer Research, Oslo University Hospital Radiumhospitalet, Oslo, Norway. <sup>100</sup>Institute of Clinical Medicine, Faculty of Medicine, University of Oslo, Oslo, Norway. <sup>101</sup>Department of Clinical Molecular Biology, Oslo University Hospital, University of Oslo, Oslo, Norway. <sup>102</sup>VIB KULeuven Center for Cancer Biology, VIB, Leuven, Belgium. <sup>103</sup>Laboratory for Translational Genetics, Department of Human Genetics, KU Leuven, Leuven, Belgium. <sup>104</sup>Epidemiology Program, University of Hawaii Cancer Center, Honolulu, HI, USA. <sup>105</sup>Department of Epidemiology, University of Washington School of Public Health, Seattle, WA, USA. <sup>106</sup>Department of Cancer Epidemiology and Prevention, M. Sklodowska-Curie Institute - Oncology Center, Warsaw, Poland. <sup>107</sup>German Breast Group, GmbH, Neu Isenburg, Germany. <sup>108</sup>Southampton Clinical Trials Unit, University of Southampton, Southampton, UK. <sup>109</sup>Research Centre for Genetic Engineering and Biotechnology "Georgi D. Efremov", Macedonian Academy of Sciences and Arts, Skopje, Republic of Macedonia. <sup>110</sup>Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>111</sup>Department of Oncology - Pathology, Karolinska Institutet, Stockholm, Sweden. <sup>112</sup>Department of Medical Oncology, University Hospital of Heraklion, Heraklion, Greece. <sup>113</sup>Division of Gynaecology and Obstetrics, Technische Universität München, Munich, Germany. <sup>114</sup>Gynaecological Cancer Research Centre, Women's Cancer, Institute for Women's Health, University College London, London, UK. <sup>115</sup>Department of Laboratory Medicine and

Pathobiology, University of Toronto, Toronto, ON, Canada. <sup>116</sup>Laboratory Medicine Program, University Health Network, Toronto, ON, Canada. <sup>117</sup>Department of Population Sciences, Beckman Research Institute of City of Hope, Duarte, CA, USA. <sup>118</sup>Department of Obstetrics and Gynecology, Helsinki University Hospital, University of Helsinki, Helsinki, Finland. <sup>119</sup>Leuven Multidisciplinary Breast Center, Department of Oncology, Leuven Cancer Institute, University Hospitals Leuven, Leuven, Belgium. <sup>120</sup>Center for Clinical Cancer Genetics and Global Health, The University of Chicago, Chicago, IL, USA. <sup>121</sup>IFOM, The FIRC (Italian Foundation for Cancer Research) Institute of Molecular Oncology, Milan, Italy. <sup>122</sup>Laboratory of Cancer Genetics and Tumor Biology, Cancer and Translational Medicine Research Unit, Biocenter Oulu, University of Oulu, Oulu, Finland. <sup>123</sup>Laboratory of Cancer Genetics and Tumor Biology, Northern Finland Laboratory Centre Oulu, Oulu, Finland. <sup>124</sup>Department of Gynecology and Obstetrics, Ludwig-Maximilians University of Munich, Munich, Germany. <sup>125</sup>Unit of Molecular Bases of Genetic Risk and Genetic Testing, Department of Research, Fondazione IRCCS (Istituto Di Ricovero e Cura a Carattere Scientifico) Istituto Nazionale dei Tumori (INT), Milan, Italy. <sup>126</sup>Section of Cancer Genetics, The Institute of Cancer Research, London, UK. <sup>127</sup>Department of Community Medicine and Epidemiology, Carmel Medical Center, Ruth and Bruce Rappaport Faculty of Medicine, Technion-Israel Institute of Technology and Clalit National Cancer Control Center, Haifa, Israel. <sup>128</sup>Medical Oncology Department, Hospital Universitario Puerta de Hierro, Madrid, Spain. <sup>129</sup>Hereditary Cancer Clinic, University Hospital of Heraklion, Heraklion, Greece. <sup>130</sup>Epidemiology Branch, National Institute of Environmental Health Sciences, NIH, Research Triangle Park, NC, USA. <sup>131</sup>Research Oncology, Guy's Hospital, King's College London, London, UK. <sup>132</sup>Division of Psychosocial Research and Epidemiology, The Netherlands Cancer Institute - Antoni van Leeuwenhoek hospital, Amsterdam, The Netherlands. <sup>133</sup>National Center for Tumor Diseases, University of Heidelberg, Heidelberg, Germany. <sup>134</sup>Division of Molecular Medicine, Pathology North, John Hunter Hospital, Newcastle, Australia. <sup>135</sup>Discipline of Medical Genetics, School of Biomedical Sciences and Pharmacy, Faculty of Health, University of Newcastle, Callaghan, Australia. <sup>136</sup>Department of Pathology, The University of Melbourne, Melbourne, Australia. <sup>137</sup>Cancer Control Research, BC Cancer Agency, Vancouver, BC, Canada. <sup>138</sup>School of Population and Public Health, University of British Columbia, Vancouver, BC, Canada. <sup>139</sup>The Curtin UWA Centre for Genetic Origins of Health and Disease, Curtin University and University of Western Australia, Perth, Australia. <sup>140</sup>Department of Obstetrics and Gynaecology, University of Melbourne and the Royal Women's Hospital, Melbourne, Australia. <sup>141</sup>Division of Breast Cancer Research, The Institute of Cancer Research, London, UK. <sup>142</sup>Epigenetic and Stem Cell Biology Laboratory, National Institute of Environmental Health Sciences, NIH, Research Triangle Park, NC, USA. <sup>143</sup>Department of Epidemiology, Mailman School of Public Health, Columbia University, New York, NY, USA. <sup>144</sup>McGill University and Génome Québec Innovation Centre, Montréal, QC, Canada. <sup>145</sup>Department of Surgery, Leiden University Medical Center, Leiden, The Netherlands. <sup>146</sup>Institute of Human

Genetics, Pontificia Universidad Javeriana, Bogota, Colombia. <sup>147</sup>Department of Gynecology and Obstetrics, Helios Clinics Berlin-Buch, Berlin, Germany. <sup>148</sup>Biostatistics and Computational Biology Branch, National Institute of Environmental Health Sciences, NIH, Research Triangle Park, NC, USA. <sup>149</sup>Department of Medicine, Institute for Human Genetics, UCSF Helen Diller Family Comprehensive Cancer Center, University of California San Francisco, San Francisco, CA, USA. <sup>150</sup>Peter MacCallum Cancer Center, Melbourne, Australia. <sup>151</sup>Department of Molecular Physiology & Biophysics, Vanderbilt Genetics Institute, Vanderbilt University, Nashville, TN, USA. <sup>152</sup>Department of Clinical Genetics, VU University Medical Center, Amsterdam, The Netherlands. <sup>153</sup>Department of Medical Oncology, Dana Farber Cancer Institute, Boston, MA. <sup>154</sup>Department of Medicine, Harvard Medical School, Boston, MA. <sup>155</sup>Division of Genetics, Brigham and Women's Hospital, Boston, MA. <sup>156</sup>Department of Nutrition, Harvard T.H. Chan School of Public Health, Boston, MA. <sup>157</sup>Human Genetics, Genome Institute of Singapore, Singapore, Singapore. <sup>158</sup>Nuffield Department of Population Health, University of Oxford, Big Data Institute, Old Road Campus, Oxford OX3 7LF, UK. <sup>159</sup>Department of Oncology University of Örebro, Örebro, Sweden. <sup>160</sup>Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>161</sup>Department of Oncology, Södersjukhuset, Stockholm, Sweden. <sup>162</sup>Independent Laboratory of Molecular Biology and Genetic Diagnostics, Pomeranian Medical University, Szczecin, Poland. <sup>163</sup>Lang Wu and Wei Shi are joint co-first authors

## Acknowledgements

The authors thank Jing He, Wanqing Wen, Ayush Giri, and Todd Edwards of Vanderbilt Epidemiology Center and Rao Tao of Department of Biostatistics, Vanderbilt University Medical Center for their help with the data analysis of this study. The authors also would like to thank all the individuals for their participation in the parent studies and all the researchers, clinicians, technicians and administrative staff for their contribution to the studies. We are also grateful to Hae Kyung Im of University of Chicago for her help. The data analyses were conducted using the Advanced Computing Center for Research and Education (ACCRE) at Vanderbilt University. This project at Vanderbilt University Medical Center was supported in part by grants R01CA158473 and R01CA148677 from the U.S. National Institutes of Health as well as funds from Anne Potter Wilson endowment. Lang Wu is supported by NCI K99 CA218892 and the Vanderbilt Molecular and Genetic Epidemiology of Cancer (MAGEC) training program (U.S. NCI grant R25 CA160056 awarded to X.-O. Shu). Genotyping of the OncoArray was principally funded from three sources: the PERSPECTIVE project, funded from the Government of Canada through Genome Canada and the Canadian Institutes of Health Research, the *Ministère de l'Économie, de la Science et de l'Innovation du Québec* through Genome Québec, and the Quebec Breast Cancer Foundation; the NCI Genetic Associations and Mechanisms in Oncology (GAME-ON) initiative and Discovery, Biology and Risk of Inherited Variants in Breast Cancer (DRIVE) project (NIH Grants U19 CA148065 and X01HG007492); and Cancer Research UK (C1287/A10118 and C1287/A16563). **BCAC** is funded by Cancer Research UK [C1287/A16563], by the European Community's Seventh Framework Programme under grant agreement 223175 (HEALTH-F2-2009-223175) (COGS) and by the European Union's Horizon 2020 Research and Innovation Programme under grant agreements 633784 (B-CAST) and 634935 (BRIDGES). Genotyping of the iCOGS array was funded by the European Union (HEALTH-F2-2009-223175), Cancer Research UK (C1287/A10710), the Canadian Institutes of Health Research for the "CIHR Team in Familial Risks of Breast Cancer" program, and the Ministry of Economic Development, Innovation and Export Trade of Quebec – grant # PSR-SIIRI-701. Combining the GWAS data was supported in part by The National Institute of Health (NIH) Cancer Post-Cancer GWAS initiative grant U19 CA 148065 (DRIVE, part of the GAME-ON initiative). A full description of funding and acknowledgments for BCAC studies, along with consortium membership are included in the Supplementary Note of the Supplementary Material.

## References

1. Kamangar F, Dores GM & Anderson WF Patterns of cancer incidence, mortality, and prevalence across five continents: defining priorities to reduce cancer disparities in different geographic regions of the world. *J Clin Oncol* 24, 2137–50 (2006). [PubMed: 16682732]
2. Beggs AD & Hodgson SV Genomics and breast cancer: the different levels of inherited susceptibility. *Eur J Hum Genet* 17, 855–6 (2009). [PubMed: 19092772]
3. Southey MC et al. PALB2, CHEK2 and ATM rare variants and cancer risk: data from COGS. *J Med Genet* (2016).
4. Nathanson KL, Wooster R & Weber BL Breast cancer genetics: what we know and what we need. *Nat Med* 7, 552–6 (2001). [PubMed: 11329055]
5. Prevalence and penetrance of BRCA1 and BRCA2 mutations in a population-based series of breast cancer cases. Anglian Breast Cancer Study Group. *Br J Cancer* 83, 1301–8 (2000).
6. Milne RL et al. Identification of ten variants associated with risk of estrogen-receptor-negative breast cancer. *Nat Genet* 49, 1767–1778 (2017). [PubMed: 29058716]
7. Michailidou K et al. Association analysis identifies 65 new breast cancer risk loci. *Nature* 551, 92–94 (2017). [PubMed: 29059683]
8. Michailidou K et al. Large-scale genotyping identifies 41 new loci associated with breast cancer risk. *Nat Genet* 45, 353–61, 361e1–2 (2013). [PubMed: 23535729]
9. Michailidou K et al. Genome-wide association analysis of more than 120,000 individuals identifies 15 new susceptibility loci for breast cancer. *Nat Genet* 47, 373–80 (2015). [PubMed: 25751625]
10. Cai Q et al. Genome-wide association analysis in East Asians identifies breast cancer susceptibility loci at 1q32.1, 5q14.3 and 15q26.1. *Nat Genet* 46, 886–90 (2014). [PubMed: 25038754]
11. Zheng W et al. Common genetic determinants of breast-cancer risk in East Asian women: a collaborative study of 23 637 breast cancer cases and 25 579 controls. *Hum Mol Genet* 22, 2539–50 (2013). [PubMed: 23535825]
12. Zhang B, Beeghly-Fadiel A, Long J & Zheng W Genetic variants associated with breast-cancer risk: comprehensive research synopsis, meta-analysis, and epidemiological evidence. *Lancet Oncol* 12, 477–88 (2011). [PubMed: 21514219]
13. French JD et al. Functional variants at the 11q13 risk locus for breast cancer regulate cyclin D1 expression through long-range enhancers. *Am J Hum Genet* 92, 489–503 (2013). [PubMed: 23540573]
14. Hindorf LA et al. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc Natl Acad Sci U S A* 106, 9362–7 (2009). [PubMed: 19474294]
15. Consortium EP An integrated encyclopedia of DNA elements in the human genome. *Nature* 489, 57–74 (2012). [PubMed: 22955616]
16. Roadmap Epigenomics C et al. Integrative analysis of 111 reference human epigenomes. *Nature* 518, 317–30 (2015). [PubMed: 25693563]
17. Dunning AM et al. Breast cancer risk variants at 6q25 display different phenotype associations and regulate ESR1, RMND1 and CCDC170. *Nat Genet* 48, 374–86 (2016). [PubMed: 26928228]
18. Ghossaini M et al. Evidence that breast cancer risk at the 2q35 locus is mediated through IGFBP5 regulation. *Nat Commun* 4, 4999 (2014). [PubMed: 25248036]
19. Li Q et al. Integrative eQTL-based analyses reveal the biology of breast cancer risk loci. *Cell* 152, 633–41 (2013). [PubMed: 23374354]
20. Darabi H et al. Polymorphisms in a Putative Enhancer at the 10q21.2 Breast Cancer Risk Locus Regulate NRBF2 Expression. *Am J Hum Genet* 97, 22–34 (2015). [PubMed: 26073781]
21. Glubb DM et al. Fine-scale mapping of the 5q11.2 breast cancer locus reveals at least three independent risk variants regulating MAP3K1. *Am J Hum Genet* 96, 5–20 (2015). [PubMed: 25529635]
22. Lawrenson K et al. Functional mechanisms underlying pleiotropic risk alleles at the 19p13.1 breast-ovarian cancer susceptibility locus. *Nat Commun* 7, 12675 (2016). [PubMed: 27601076]

23. Lee D et al. A method to predict the impact of regulatory variants from DNA sequence. *Nat Genet* 47, 955–61 (2015). [PubMed: 26075791]
24. Finucane HK et al. Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat Genet* 47, 1228–35 (2015). [PubMed: 26414678]
25. Gusev A et al. Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am J Hum Genet* 95, 535–52 (2014). [PubMed: 25439723]
26. Barbeira AN et al. Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *bioRxiv* (2017).
27. Gamazon ER et al. A gene-based association method for mapping traits using reference transcriptome data. *Nat Genet* 47, 1091–8 (2015). [PubMed: 26258848]
28. Gusev A et al. Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet* 48, 245–52 (2016). [PubMed: 26854917]
29. Zhu Z et al. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat Genet* 48, 481–7 (2016). [PubMed: 27019110]
30. Hoffman JD et al. Cis-eQTL-based trans-ethnic meta-analysis reveals novel genes associated with breast cancer risk. *PLoS Genet* 13, e1006690 (2017). [PubMed: 28362817]
31. Lin WY et al. Identification and characterization of novel associations in the CASP8/ALS2CR12 region on chromosome 2 with breast cancer risk. *Hum Mol Genet* 24, 285–98 (2015). [PubMed: 25168388]
32. Camp NJ et al. Discordant Haplotype Sequencing Identifies Functional Variants at the 2q33 Breast Cancer Risk Locus. *Cancer Res* 76, 1916–25 (2016). [PubMed: 26795348]
33. Li Q et al. Expression QTL-based analyses reveal candidate causal genes and loci across five tumor types. *Hum Mol Genet* 23, 5294–302 (2014). [PubMed: 24907074]
34. Caswell JL et al. Multiple breast cancer risk variants are associated with differential transcript isoform expression in tumors. *Hum Mol Genet* 24, 7421–31 (2015). [PubMed: 26472073]
35. Darabi H et al. Fine scale mapping of the 17q22 breast cancer locus using dense SNPs, genotyped within the Collaborative Oncological Gene-Environment Study (COGs). *Sci Rep* 6, 32512 (2016). [PubMed: 27600471]
36. Yang J et al. Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet* 44, 369–75, S1–3 (2012). [PubMed: 22426310]
37. Kramer A, Green J, Pollard J, Jr. & Tugendreich S Causal analysis approaches in Ingenuity Pathway Analysis. *Bioinformatics* 30, 523–30 (2014). [PubMed: 24336805]
38. Koh JL et al. COLT-Cancer: functional genetic screening resource for essential genes in human cancer cell lines. *Nucleic Acids Res* 40, D957–63 (2012). [PubMed: 22102578]
39. Marcotte R et al. Essential gene profiles in breast, pancreatic, and ovarian cancer cells. *Cancer Discov* 2, 172–89 (2012). [PubMed: 22585861]
40. Walen KH & Stampfer MR Chromosome analyses of human mammary epithelial cells at stages of chemical-induced transformation progression to immortality. *Cancer Genet Cytogenet* 37, 249–61 (1989). [PubMed: 2702624]
41. Tszetzamsky AD et al. BRCA1- and BRCA2-deficient cells are sensitive to etoposide-induced DNA double-strand breaks via topoisomerase II. *Cancer Res* 67, 7078–81 (2007). [PubMed: 17671173]
42. Marcotte R et al. Essential gene profiles in breast, pancreatic, and ovarian cancer cells. *Cancer Discov* 2, 172–189 (2012). [PubMed: 22585861]
43. Sanchez Y et al. Genome-wide analysis of the human p53 transcriptional network unveils a lncRNA tumour suppressor signature. *Nat Commun* 5, 5812 (2014). [PubMed: 25524025]
44. Li Y, Peart MJ & Prives C Stxbp4 regulates DeltaNp63 stability by suppression of RACK1-dependent degradation. *Mol Cell Biol* 29, 3953–63 (2009). [PubMed: 19451233]
45. Sekine Y et al. The Kelch repeat protein KLHDC10 regulates oxidative stress-induced ASK1 activation by suppressing PP5. *Mol Cell* 48, 692–704 (2012). [PubMed: 23102700]

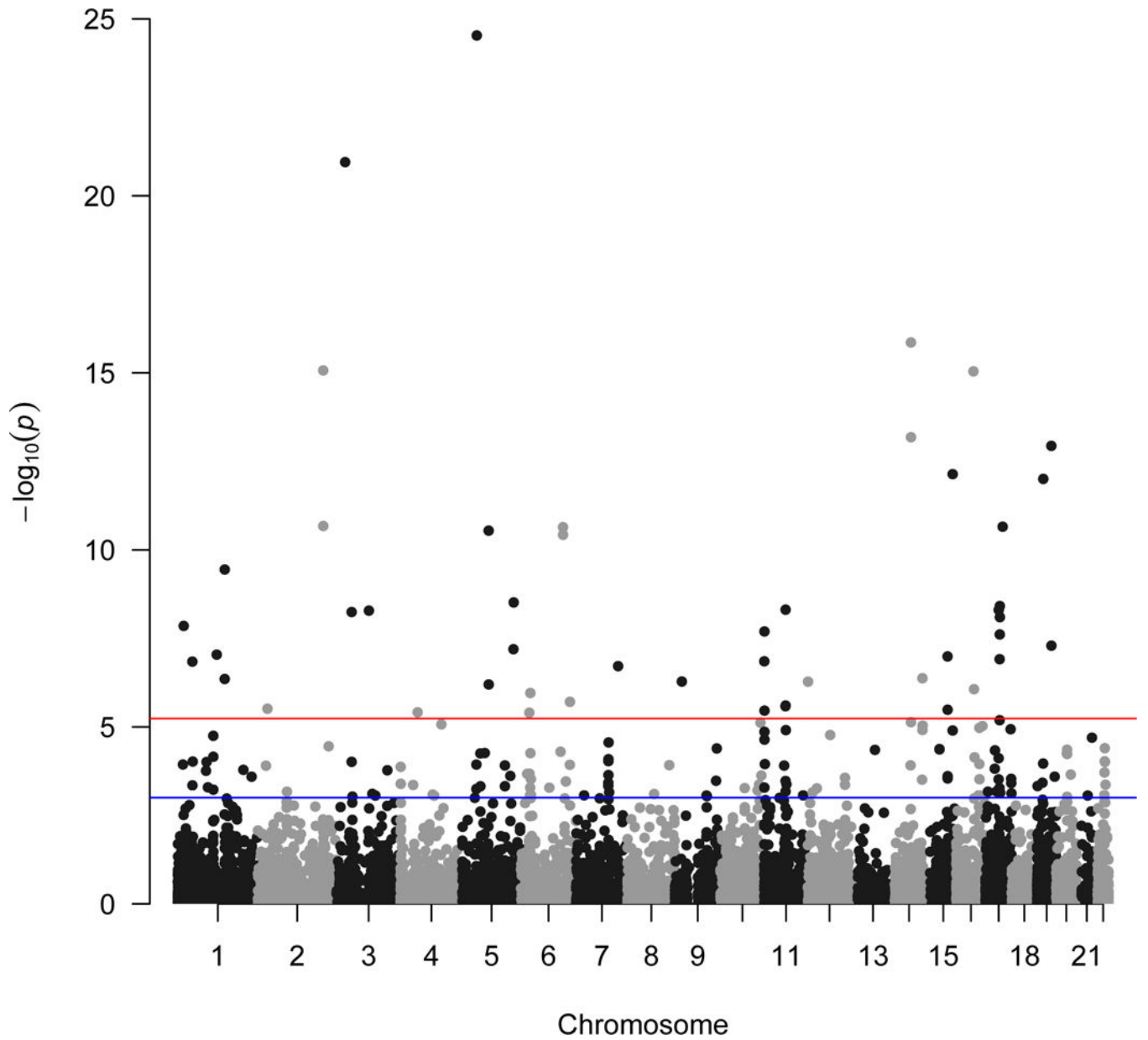


46. Kim MH et al. Anaplastic lymphoma kinase gene copy number gain in inflammatory breast cancer (IBC): prevalence, clinicopathologic features and prognostic implication. *PLoS One* 10, e0120320 (2015). [PubMed: 25803816]
47. Crizotinib versus Chemotherapy in Advanced ALK-Positive Lung Cancer. *N Engl J Med* 373, 1582 (2015).
48. Le Page C et al. *BTN3A2* expression in epithelial ovarian cancer is associated with higher tumor infiltrating T cells and a better prognosis. *PLoS One* 7, e38541 (2012). [PubMed: 22685580]
49. Kan L et al. *LRRC3B* is downregulated in non-small-cell lung cancer and inhibits cancer cell proliferation and invasion. *Tumour Biol* 37, 1113–20 (2016). [PubMed: 26276358]
50. Cox A et al. A common coding variant in *CASP8* is associated with breast cancer risk. *Nat Genet* 39, 352–8 (2007). [PubMed: 17293864]
51. Yang J et al. Genomic inflation factors under polygenic inheritance. *Eur J Hum Genet* 19, 807–12 (2011). [PubMed: 21407268]
52. Marouli E et al. Rare and low-frequency coding variants alter human adult height. *Nature* 542, 186–190 (2017). [PubMed: 28146470]
53. Turcot V et al. Protein-altering variants associated with body mass index implicate pathways that control energy intake and expenditure in obesity. *Nat Genet* 50, 26–41 (2018). [PubMed: 29273807]
54. Mele M et al. Human genomics. The human transcriptome across tissues and individuals. *Science* 348, 660–5 (2015). [PubMed: 25954002]

## References

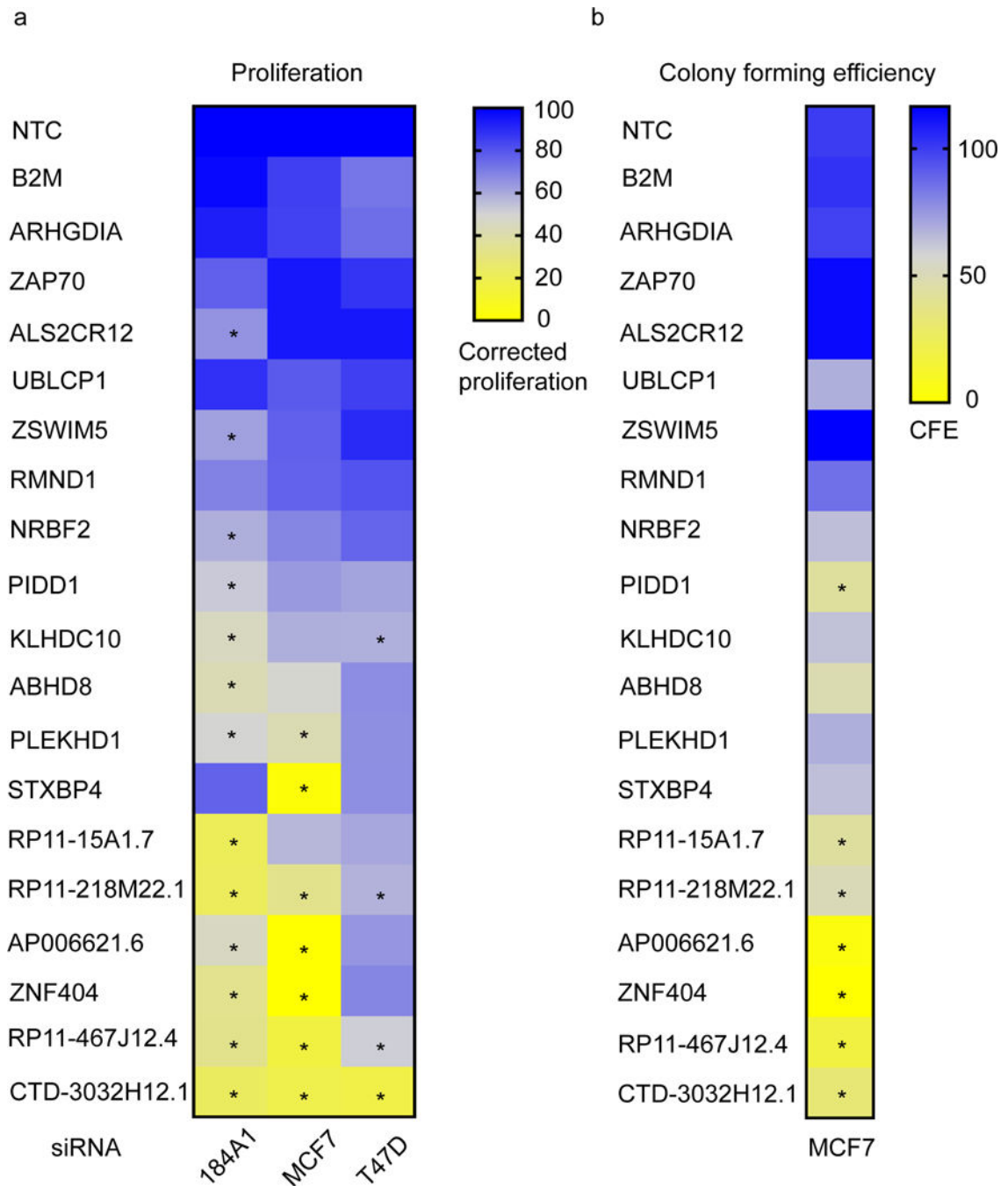
55. Consortium GT Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* 348, 648–60 (2015). [PubMed: 25954001]
56. McCarthy S et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* 48, 1279–83 (2016). [PubMed: 27548312]
57. Delaneau O, Marchini J & Zagury JF A linear complexity phasing method for thousands of genomes. *Nat Methods* 9, 179–81 (2012).
58. Howie BN, Donnelly P & Marchini J A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet* 5, e1000529 (2009). [PubMed: 19543373]
59. DeLuca DS et al. RNA-SeQC: RNA-seq metrics for quality control and process optimization. *Bioinformatics* 28, 1530–2 (2012). [PubMed: 22539670]
60. Stegle O, Parts L, Piipari M, Winn J & Durbin R Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat Protoc* 7, 500–7 (2012). [PubMed: 22343431]
61. Guo X, Lin M, Rockowitz S, Lachman HM & Zheng D Characterization of human pseudogene-derived non-coding RNAs for functional potential. *PLoS One* 9, e93972 (2014). [PubMed: 24699680]
62. Casbas-Hernandez P et al. Tumor intrinsic subtype is reflected in cancer-adjacent tissue. *Cancer Epidemiol Biomarkers Prev* 24, 406–14 (2015). [PubMed: 25465802]
63. Huang X, Stern DF & Zhao H Transcriptional Profiles from Paired Normal Samples Offer Complementary Information on Cancer Patient Survival--Evidence from TCGA Pan-Cancer Data. *Sci Rep* 6, 20567 (2016). [PubMed: 26837275]
64. Ghousaini M et al. Genome-wide association analysis identifies three new breast cancer susceptibility loci. *Nat Genet* 44, 312–8 (2012). [PubMed: 22267197]
65. Garcia-Closas M et al. Genome-wide association studies identify four ER negative-specific breast cancer risk loci. *Nat Genet* 45, 392–8, 398e1–2 (2013). [PubMed: 23535733]
66. Devlin B & Roeder K Genomic control for association studies. *Biometrics* 55, 997–1004 (1999). [PubMed: 11315092]
67. Freedman ML et al. Assessing the impact of population stratification on genetic association studies. *Nat Genet* 36, 388–93 (2004). [PubMed: 15052270]

68. Rao SS et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* 159, 1665–80 (2014). [PubMed: 25497547]
69. He B, Chen C, Teng L & Tan K Global view of enhancer-promoter interactome in human cells. *Proc Natl Acad Sci U S A* 111, E2191–9 (2014). [PubMed: 24821768]
70. Corradin O et al. Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res* 24, 1–13 (2014). [PubMed: 24196873]
71. Hnisz D et al. Super-enhancers in the control of cell identity and disease. *Cell* 155, 934–47 (2013). [PubMed: 24119843]
72. Consortium F et al. A promoter-level mammalian expression atlas. *Nature* 507, 462–70 (2014). [PubMed: 24670764]



**Figure 1. Manhattan plot of association results from the breast cancer transcriptome-wide association study.**

Results are based on 122,977 cases and 105,974 controls. The red line represents  $P = 5.82 \times 10^{-6}$ . The blue line represents  $P = 1.00 \times 10^{-3}$ .



**Figure 2. Heat maps of proliferation and colony formation efficiency in breast cells.** (a) Proliferation efficiency. (b) colony formation efficiency. Error bars, SD ( $N=2$ ).  $P$ -values were determined by one-way ANOVA followed by Dunnett's multiple comparisons test: \* $P$ -value < 0.05. NTC: non-target control.

Table 1.

Fourteen expression-trait associations for genes located at genomic loci at least 500 kb away from any GWAS-identified breast cancer risk variants

Region	Gene <sup>a</sup>	Type <sup>b</sup>	Z score	P value <sup>c</sup>	R <sup>2</sup>	Closest risk SNP <sup>d</sup>	Distance to the closest risk SNP (kb)	P value after adjusting for adjacent risk SNPs <sup>e</sup>
1p34.1	<i>ZSWIM5</i>	Protein	5.26	$1.43 \times 10^{-7}$	0.17	rs1707302	829	0.006
3p24.1	<i>LRR3B</i>	Protein	-9.57	$1.11 \times 10^{-21}$	0.17	rs653465	591	$1.60 \times 10^{-6}$
4q12	<i>SPATA18</i>	Protein	-4.62	$3.86 \times 10^{-6}$	0.11	rs6815814	14,101	$3.98 \times 10^{-6}$
6p22.1	<i>UBD</i>	Protein	-4.87	$1.10 \times 10^{-6}$	0.13	rs9257408	597	0.94
7q32.2	<i>KLHDC10</i>	Protein	5.21	$1.92 \times 10^{-7}$	0.14	rs4593472	892	$2.90 \times 10^{-7}$
9p21.3	<i>MIR31HG</i>	lncRNA	-5.02	$5.22 \times 10^{-7}$	0.12	rs1011970	502	$1.23 \times 10^{-7}$
11p15.5	<i>RIC8A</i>	Protein	-5.27	$1.40 \times 10^{-7}$	0.15	rs6597981	588	$4.95 \times 10^{-6}$
11q13.2	<i>B3GNT1</i>	Protein	-5.85	$4.88 \times 10^{-9}$	0.09	rs3903072	530	$3.50 \times 10^{-6}$
11q13.2	<i>RP11-867G23.10</i>	transcript	4.71	$2.49 \times 10^{-6}$	0.03	rs3903072	594	$2.61 \times 10^{-4}$
12p13.33	<i>RP11-218M22.1</i>	lncRNA	5.02	$5.27 \times 10^{-7}$	0.19	rs12422552	13,641	$5.17 \times 10^{-7}$
14q24.1	<i>GALNT16</i>	Protein	-8.27	$1.38 \times 10^{-16}$	0.04	rs999737	691	$8.57 \times 10^{-4}$
14q24.1	<i>PLEKHDI</i>	Protein	7.50	$6.55 \times 10^{-14}$	0.02	rs999737	917	0.12
15q24.2	<i>MAN2C1</i> <sup>f</sup>	Protein	-5.32	$1.02 \times 10^{-7}$	0.39	rs2290203	15,851	$9.56 \times 10^{-8}$
15q24.2	<i>CTD-232K18.1</i> <sup>f</sup>	lncRNA	-4.65	$3.27 \times 10^{-6}$	0.07	rs2290203	15,619	$3.16 \times 10^{-6}$

<sup>a</sup>Genes that were siRNA-silenced for functional assays are bolded; SNPs used to predict gene expression are listed in the Supplementary Table 13

<sup>b</sup>Protein: protein coding genes; lncRNA: long non-coding RNAs; transcript: processed transcript

<sup>c</sup>P-value: derived from association analyses of 122,977 cases and 105,974 controls; associations with  $p < 5.82 \times 10^{-6}$  considered statistically significant based on Bonferroni correction of 8,597 tests (0.05/8,597); R<sup>2</sup>: prediction performance (R<sup>2</sup>) derived using GTEX data.

<sup>d</sup>Risk SNPs identified in previous GWAS or fine-mapping studies. The risk SNP closest to the gene is presented. A full list of all risk SNPs, and their distances to the genes are presented in the Supplementary Table 4

<sup>e</sup>Use of COIO method<sup>36</sup>

<sup>f</sup>Predicted expression of *MAN2C1* and *CTD-232K18.1* was correlated (spearman R=0.76)

Table 2.

Twenty-three expression-trait associations for genes located at genomic loci within 500 kb of any previous GWAS-identified breast cancer risk variants but not yet implicated as target genes of risk variants<sup>#</sup>

Region	Gene <sup>a</sup>	Type <sup>b</sup>	Z score	P value <sup>c</sup>	R <sup>2c</sup>	Closest risk SNP <sup>d</sup>	Distance to the closest risk SNP (kb)	P value after adjusting for adjacent risk SNPs <sup>e</sup>
1p11.2	<i>RPI1-439A17.7</i>	lncRNA	-5.34	$9.07 \times 10^{-8}$	0.22	rs11249433	442	0.02
1q21.1	<i>NUDT17</i>	Protein	-6.27	$3.58 \times 10^{-10}$	0.01	rs12405132	56	0.08
1q21.1	<i>ANKRD34A</i>	Protein	-5.05	$4.42 \times 10^{-7}$	0.01	rs12405132	169	$4.28 \times 10^{-5}$
2p23.1-2p23.2	<i>ALK</i>	Protein	4.67	$3.06 \times 10^{-6}$	0.06	rs4577244	295	$2.70 \times 10^{-6}$
3p21.31	<i>PRSS46</i>	Protein	-5.83	$5.68 \times 10^{-9}$	0.13	rs6796502	89	0.002
3q12.2	<i>RPI1-114f8.4</i>	lncRNA	-5.84	$5.19 \times 10^{-9}$	0.02	rs9833888	356	0.09
5p12	<i>RPI1-53O19.1</i>	lncRNA	10.38	$2.94 \times 10^{-25}$	0.03	rs10941679	39	$7.46 \times 10^{-4}$
5q33.3	<i>UBLCPI</i>	Protein	5.93	$3.04 \times 10^{-9}$	0.07	rs1432679	446	0.37
5q33.3	<i>RPI1-32D16.1</i>	lncRNA	-5.41	$6.37 \times 10^{-8}$	0.09	rs1432679	283	$1.32 \times 10^{-4}$
6p22.2	<i>BTN3A2</i>	Protein	4.61	$3.97 \times 10^{-6}$	0.28	rs71557345	229	0.72
6q23.1	<i>RPI1-73O6.3<sup>f</sup></i>	lncRNA	-6.61	$3.74 \times 10^{-11}$	0.11	rs6569648	105	0.41
11p15.5	<i>AP006621.6<sup>g</sup></i>	lncRNA	5.61	$2.01 \times 10^{-8}$	0.34	rs6597981	21	0.52
11p15.5	<i>RPLP2<sup>g</sup></i>	Protein	4.64	$3.46 \times 10^{-6}$	0.27	rs6597981	7	0.51
14q32.33	<i>CTD-3051D23.1</i>	lncRNA	-5.06	$4.21 \times 10^{-7}$	0.05	rs10623258	97	$7.05 \times 10^{-7}$
16q12.2	<i>RPI1-467J12.4</i>	lncRNA	8.04	$9.02 \times 10^{-16}$	0.23	rs3112612	434	0.79
16q12.2	<i>CTD-3032H12.1</i>	lncRNA	4.92	$8.58 \times 10^{-7}$	0.03	rs28539243	290	0.006
17q21.31	<i>LRRCS7A<sup>g</sup></i>	Protein	-5.89	$3.85 \times 10^{-9}$	0.43	rs2532263	118	0.79
17q21.31	<i>KANSL1-AS1<sup>g</sup></i>	lncRNA	-5.58	$2.44 \times 10^{-8}$	0.62	rs2532263	18	0.95
17q21.31	<i>CRHR1<sup>g</sup></i>	Protein	-5.29	$1.22 \times 10^{-7}$	0.22	rs2532263	339	0.99
17q21.31	<i>LINC00671</i>	lncRNA	-5.85	$4.95 \times 10^{-9}$	0.07	rs72826962	190	0.26
17q21.31	<i>LRRCS7A2</i>	Protein	-5.77	$7.93 \times 10^{-9}$	0.46	rs2532263	336	0.93
19p13.11	<i>HAPLN4</i>	Protein	-7.13	$9.88 \times 10^{-13}$	0.02	rs2965183	172	0.22

Region	Gene <sup>a</sup>	Type <sup>b</sup>	Z score	P value <sup>c</sup>	R <sup>2c</sup>	Closest risk SNP <sup>d</sup>	Distance to the closest risk SNP (kb)	P value after adjusting for adjacent risk SNPs <sup>e</sup>
19q13.31	<b><i>RPI1-15A1.7</i></b> <sup>h</sup>	lncRNA	5.45	5.06 × 10 <sup>-8</sup>	0.02	rs3760982	215	0.28

# not yet reported from eQTL and/or functional studies as target genes of GWAS-identified risk variants and not harbor GWAS or fine-mapping identified risk variants

<sup>a</sup> Genes that were siRNA-silenced for functional assays are bolded; SNPs used to predict gene expression are listed in the Supplementary Table 13

<sup>b</sup> Protein: protein coding genes; lncRNA: long non-coding RNAs

<sup>c</sup> P value: nominal P value from association analysis of 122,977 cases and 105,974 controls; the threshold after Bonferroni correction of 8,597 tests (0.05/8,597=5.82×10<sup>-6</sup>) was used; R<sup>2</sup>: prediction performance (R<sup>2</sup>) derived using GTEx data

<sup>d</sup> Risk SNPs identified in previous GWAS or fine-mapping studies. The risk SNP closest to the gene is presented. A full list of all risk SNPs, and their distances to the genes are presented in the Supplementary Table 4

<sup>e</sup> Use of COJO method<sup>36</sup>; all index SNPs in the corresponding region were adjusted in the conditional analyses

<sup>f</sup> Predicted expression of *RPI1-7306.3* and *L3MBTL3* was correlated (spearman R=0.88)

<sup>g</sup> Predicted expression of *AP006621.6* and *RPLP2* was correlated; predicted expression of *LRRC37A*, *KANSLL-AS1*, and *CRHR1* was correlated (spearman R>0.1)

<sup>h</sup> Predicted expression of *RPI1-15A1.7* and *ZNF404* was correlated (spearman R=0.64)

Table 3.

Eleven expression-trait associations for genes previously reported as potential target genes of GWAS-identified breast cancer risk variants or genes harboring risk variants

Region	Gene <sup>a</sup>	Type <sup>b</sup>	Z score	P value <sup>c</sup>	R <sup>2</sup> <sup>c</sup>	Closest risk SNP <sup>d</sup>	Distance to the closest risk SNP (kb)	P value after adjusting for adjacent risk SNPs <sup>e</sup>	Association direction reported previously <sup>f</sup>	Reference
1p36.13	<i>KLHDC7A</i>	Protein	-5.67	$1.40 \times 10^{-8}$	0.04	rs2992756	0.085	0.06	-	7
2q33.1	<i>ALS2CR12</i>	Protein	6.70	$2.11 \times 10^{-11}$	0.10	rs1830298	intron of the gene	0.17	NA	31
2q33.1	<i>CASP8</i>	Protein	-8.05	$8.51 \times 10^{-16}$	0.22	rs3769821	intron of the gene	0.16	-	31,32
5q14.1	<i>ATG10</i>	Protein	-6.65	$2.85 \times 10^{-11}$	0.51	rs7707921	intron of the gene	0.21	NA	9
5q14.2	<i>ATP6AP1L</i>	Protein	-4.98	$6.32 \times 10^{-7}$	0.63	rs7707921	37	0.98	NA	9
6q23.1	<i>L3MBTL3</i> <sup>g</sup>	Protein	-6.69	$2.27 \times 10^{-11}$	0.10	rs6569648	208	0.44	NA	6
6q25.1	<i>RMND1</i>	Protein	4.76	$1.95 \times 10^{-6}$	0.13	rs3757322	169	$1.11 \times 10^{-4}$	mixed	17
11q13.1	<i>SNX32</i>	Protein	4.70	$2.60 \times 10^{-6}$	0.19	rs3903072	18	0.17	NA	33
15q26.1	<i>RCCDI</i>	Protein	-7.18	$7.23 \times 10^{-13}$	0.13	rs2290203	6	$1.66 \times 10^{-4}$	-	10
17q22	<i>STXBPA</i>	Protein	6.69	$2.21 \times 10^{-11}$	0.03	rs6504950	intron of the gene	0.90	+ in GTEx	34,35
19q13.31	<i>ZNF404</i> <sup>h</sup>	Protein	7.42	$1.15 \times 10^{-13}$	0.15	rs3760982	90	0.005	NA	8

<sup>a</sup> Genes that were siRNA silenced for functional assays are bolded; SNPs used to predict gene expression are listed in the Supplementary Table I3

<sup>b</sup> Protein: protein coding genes; lncRNA: long non-coding RNAs; NA: not available

<sup>c</sup> P value: nominal P value from association analysis of 122,977 cases and 105,974 controls; the threshold after Bonferroni correction of 8,597 tests ( $0.05/8,597=5.82 \times 10^{-6}$ ) was used; R<sup>2</sup>: prediction performance (R<sup>2</sup>) derived using GTEx data.

<sup>d</sup> Risk SNPs identified in previous GWAS or fine-mapping studies. The risk SNP closest to the gene is presented. A full list of all risk SNPs, and their distances to the genes are presented in the Supplementary Table 4

<sup>e</sup> Use of COIO method<sup>36</sup>; all index SNPs in the corresponding region were adjusted for the conditional analyses

<sup>f</sup> -: inverse association; +: positive association; mixed: both inverse and positive associations reported; NA: not available

<sup>g</sup> Predicted expression of *L3MBTL3* and *RPI1-7306.3* was correlated (spearman R=0.88)

<sup>h</sup> Predicted expression of *ZNF404* and *RPI1-15A1.7* was correlated (spearman R=0.64)



**Table 4.**

Genes at GWAS-identified breast cancer risk loci ( $\pm 500$ kb of the index SNPs) whose predicted expression levels were associated with breast cancer risk at  $p$ -values between  $5.82 \times 10^{-6}$  and  $1.05 \times 10^{-3}$  (FDR corrected  $p$ -value 0.05)

Region	Gene	Type <sup>a</sup>	Z score	$P$ value <sup>b</sup>	$R^2$ <sup>b</sup>	Closest risk SNP <sup>c</sup>	Distance to the closest risk SNP (kb)	$P$ value after adjusting for adjacent risk SNPs <sup>d</sup>
1p34.1	<i>UQCRH</i>	Protein	-3.90	$9.51 \times 10^{-5}$	0.12	rs1707302	168	0.06
1p22.3	<i>LMO4</i>	Protein	-3.76	$1.73 \times 10^{-4}$	0.09	rs12118297	15	0.002
2p23.3	<i>DNAJC27-AS1</i>	lncRNA	3.84	$1.24 \times 10^{-4}$	0.03	rs6725517	65	0.13
4p14	<i>KLHL5</i>	Protein	3.52	$4.35 \times 10^{-4}$	0.13	rs6815814	230	0.03
5q11.2	<i>AC008391.1</i>	miRNA	-4.03	$5.60 \times 10^{-5}$	0.13	rs16886113	242	0.76
6p22.1	<i>HCG14</i>	lncRNA	-3.47	$5.19 \times 10^{-4}$	0.11	rs9257408	61	0.03
6p22.2	<i>TRNAI2</i>	miRNA	-3.71	$2.09 \times 10^{-4}$	0.02	rs71557345	307	0.007
6q25.1	<i>MTHFD1L</i>	Protein	3.85	$1.17 \times 10^{-4}$	0.10	rs3757318	491	$2.36 \times 10^{-4}$
8q24.21	<i>PVT1</i>	transcript	3.85	$1.20 \times 10^{-4}$	0.03	rs11780156	81	$1.09 \times 10^{-4}$
9q33.3	<i>RP11-123K19.1</i>	lncRNA	-4.10	$4.05 \times 10^{-5}$	0.05	rs10760444	20	$1.26 \times 10^{-4}$
10q25.2	<i>RP11-57H14.3</i>	lncRNA	3.42	$6.16 \times 10^{-4}$	0.08	rs7904519	108	0.002
10q26.13	<i>RP11-500G22.2</i>	lncRNA	4.48	$7.54 \times 10^{-6}$	0.15	rs2981582	336	0.91
11p15.5	<i>PTDSS2</i>	Protein	-3.47	$5.16 \times 10^{-4}$	0.04	rs6597981	312	0.02
11p15.5	<i>AP006621.5</i>	Protein	4.35	$1.37 \times 10^{-5}$	0.51	rs6597981	19	0.01
11p15.5	<i>PIDD1</i>	Protein	4.24	$2.28 \times 10^{-5}$	0.45	rs6597981	intron of the gene	0.12
11p15.5	<i>MRPL23-AS1</i>	lncRNA	-3.86	$1.12 \times 10^{-4}$	0.10	rs3817198	95	0.06
11q13.1-11q13.2	<i>PACSI1</i>	Protein	-3.59	$3.36 \times 10^{-4}$	0.06	rs3903072	255	0.001
12p11.22	<i>RP11-860B13.1</i>	lncRNA	3.46	$5.42 \times 10^{-4}$	0.17	rs10771399	221	0.86
13q22.1	<i>KLIF5</i>	Protein	-4.08	$4.44 \times 10^{-5}$	0.22	rs6562760	306	NA
14q24.1	<i>CTD-256613.1</i>	lncRNA	-3.84	$1.22 \times 10^{-4}$	0.04	rs2588809	64	0.55
14q32.33	<i>C14orf79</i>	Protein	4.37	$1.22 \times 10^{-5}$	0.11	rs10623258	240	0.91
15q26.1	<i>FES</i>	Protein	4.37	$1.26 \times 10^{-5}$	0.21	rs2290203	73	$3.04 \times 10^{-6}$
16q12.2	<i>BBS2</i>	Protein	3.97	$7.23 \times 10^{-5}$	0.26	rs2432539	80	0.36

Region	Gene	Type <sup>a</sup>	Z score	P value <sup>b</sup>	R <sup>2b</sup>	Closest risk SNP <sup>c</sup>	Distance to the closest risk SNP (kb)	P value after adjusting for adjacent risk SNPs <sup>d</sup>
16q12.2	<i>CRNDE</i>	lncRNA	3.28	$1.05 \times 10^{-3}$	0.02	rs28539243	271	0.69
16q24.2	<i>RP11-482M8.1</i>	lncRNA	3.32	$9.16 \times 10^{-4}$	0.02	rs4496150	441	0.19
17q11.2	<i>GOSR1</i>	Protein	3.79	$1.51 \times 10^{-4}$	0.10	rs146699004	376	0.04
17q21.2	<i>ATP6V0A1</i>	Protein	3.61	$3.02 \times 10^{-4}$	0.03	rs72826962	162	0.01
17q21.2	<i>RP11-400F19.8</i>	transcript	-3.96	$7.65 \times 10^{-5}$	0.01	rs72826962	122	$6.62 \times 10^{-4}$
17q21.31	<i>RP11-105N13.4</i>	transcript	-4.51	$6.46 \times 10^{-6}$	0.02	rs252263	359	NA
17q25.3	<i>CBX8</i>	Protein	4.38	$1.16 \times 10^{-5}$	0.05	rs745570	6	0.99
19p13.11	<i>CTD-2538G9.5</i>	lncRNA	3.56	$3.76 \times 10^{-4}$	0.01	rs8170	432	$4.38 \times 10^{-4}$
19p13.11	<i>HOMER3</i>	Protein	-3.87	$1.08 \times 10^{-4}$	0.10	rs4808801	469	0.18
20q11.22	<i>CTD-3216D2.5</i>	lncRNA	4.03	$5.60 \times 10^{-5}$	0.16	rs2284378	281	$9.24 \times 10^{-4}$
22q13.1	<i>TRIOBP</i>	Protein	3.34	$8.34 \times 10^{-4}$	0.07	rs738321	396	0.003
22q13.1	<i>RP5-1039K5.13</i>	lncRNA	3.73	$1.93 \times 10^{-4}$	0.01	rs738321	99	0.053
22q13.1	<i>CBY1</i>	Protein	3.91	$9.34 \times 10^{-5}$	0.05	chr22:39359355	289	0.06
22q13.1	<i>APOBEC3A</i>	Protein	-4.11	$3.98 \times 10^{-5}$	0.07	chr22:39359355	0.2	0.02
22q13.2	<i>RP1-85F18.6</i>	lncRNA	3.52	$4.28 \times 10^{-4}$	0.12	rs73161324	460	0.72

<sup>a</sup>Protein: protein coding genes; lncRNA: long non-coding RNAs; transcript: processed transcript

<sup>b</sup>P value: nominal P value from association analysis of 122,977 cases and 105,974 controls; R<sup>2</sup>: prediction performance derived using GTEx data.

<sup>c</sup>Risk SNPs identified in previous GWAS or fine-mapping studies. The risk SNP closest to the gene is presented. A full list of all risk SNPs, and their distances to the genes are presented in the Supplementary Table 4

<sup>d</sup>Use of COJO method<sup>36</sup>; all index SNPs in the corresponding region were adjusted for the conditional analyses